



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Accelerated Stochastic Matrix Inversion: General Theory and Speeding up BFGS Rules for Faster Second-Order Optimization

Citation for published version:

Gower, RM, Hanzely, F, Richtárik, P & Stich, S 2018, 'Accelerated Stochastic Matrix Inversion: General Theory and Speeding up BFGS Rules for Faster Second-Order Optimization', Thirty-second Conference on Neural Information Processing Systems, Montreal, Canada, 3/12/18 - 8/12/18.
<<http://papers.nips.cc/paper/7434-accelerated-stochastic-matrix-inversion-general-theory-and-speeding-up-bfgs-rules-for-faster-second-order-optimization>>

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Early version, also known as pre-print

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Accelerated Stochastic Matrix Inversion: General Theory and Speeding up BFGS Rules for Faster Second-Order Optimization

Robert M. Gower^{*} Filip Hanzely[†] Peter Richtárik[‡] Sebastian Stich[§]

February 12, 2018

Abstract

We present the first accelerated randomized algorithm for solving linear systems in Euclidean spaces. One essential problem of this type is the matrix inversion problem. In particular, our algorithm can be specialized to invert positive definite matrices in such a way that all iterates (approximate solutions) generated by the algorithm are positive definite matrices themselves. This opens the way for many applications in the field of optimization and machine learning. As an application of our general theory, we develop the *first accelerated (deterministic and stochastic) quasi-Newton updates*. Our updates lead to provably more aggressive approximations of the inverse Hessian, and lead to speed-ups over classical non-accelerated rules in numerical experiments. Experiments with empirical risk minimization show that our rules can accelerate training of machine learning models.

1 Introduction

Second order methods have played a key role in the field of optimization throughout its history. The simplest, Newton’s method

$$w_{k+1} = w_k - (\nabla^2 f(w_k))^{-1} \nabla f(w_k), \quad (1)$$

is however inefficient for solving larger problems as it requires an access to the whole Hessian and then solve linear system each iteration. Several methods have been developed to address this issue, mostly approximating the exact update above.

Quasi-Newton methods, in particular the BFGS [3, 6, 7, 25], have been the leading optimization algorithm in various fields since late 60’s until the rise of big data, which brought a need for simpler first order algorithms. It is well known that Nesterov’s acceleration [19] is a reliable way to speed up first order methods. However until now, acceleration techniques have been applied exclusively to speeding up gradient updates. In this paper we present an accelerated BFGS algorithm, opening up new applications for acceleration. The acceleration in fact comes from an accelerated algorithm for inverting the Hessian matrix.

^{*}Télécom ParisTech, Paris, France

[†]King Abdullah University of Science and Technology, Thuwal, Saudi Arabia

[‡]King Abdullah University of Science and Technology, Thuwal, Saudi Arabia — University of Edinburgh, Edinburgh, United Kingdom — Moscow Institute of Physics and Technology, Moscow, Russia

[§]École polytechnique fédérale de Lausanne (EPFL), Lausanne, Switzerland

To be more specific, recall that quasi-Newton rules aim to maintain an estimate of the inverse Hessian X_k , adjusting it every iteration so that the inverse Hessian acts appropriately in a particular direction, while enforcing symmetry:

$$\begin{aligned} X_k(w_k - w_{k-1}) &= \nabla f(w_k) - \nabla f(w_{k-1}) \\ X_k &= X_k^\top. \end{aligned} \quad (2)$$

A notable research direction is the development of stochastic quasi-Newton methods [12], where the estimated inverse is equal to the true inverse over a subspace:

$$X_k \nabla^2 f(w_k) S_k = S_k, \quad X_k = X_k^\top, \quad (3)$$

where $S_k \in \mathbb{R}^{n \times \tau}$ is a randomly generated matrix.

In fact, (3) can be seen as the so called sketch-and-project iteration for inverting $\nabla^2 f(w_k)$. In this paper we first develop the accelerated algorithm for inverting positive definite matrices. As a direct application, our algorithm can be used as a primitive in quasi-Newton methods to expedite the Newton step as we will demonstrate. This results in a novel accelerated (stochastic) quasi-Newton method of the type (3). In addition, our acceleration technique can also be incorporated in the classical (non stochastic) BFGS updates. This results in the accelerated BFGS method. Whereas the matrix inversion contribution is accompanied by strong theoretical justifications, this does not apply to the latter. Rather, we verify the effectiveness of this new accelerated BFGS method through numerical experiments.

We should also mention an increasing recent development of second order methods to solve large optimization problems. In particular, stochastic BFGS [17, 30] is one of efficient approaches. On the other hand, *sketching* approaches as in [22, 32] reduce the dimension, and hence the complexity of the Hessian and the update, whereas *subsampling* Newton methods reduce the complexity by exploiting additional structure of the Hessian (for instance when the optimization objective can be written as the sum of many loss functions) [2, 1].

1.1 Sketch-and-project for linear systems

Our accelerated algorithm can be applied to more general tasks than only inverting matrices. In its most general form, it can be seen as an accelerated version of a *sketch-and-project* method in Euclidean spaces which we present now. Consider a linear system $Ax = b$ such that $b \in \text{Range}(A)$. One step of the sketch-and-project algorithm reads as:

$$\begin{aligned} x_{k+1} &= \operatorname{argmin}_x \|x_k - x\|_B^2 \\ \text{subject to } & S_k^\top Ax = S_k^\top b, \end{aligned}$$

where $\|x\|_B = \langle Bx, x \rangle$ and S_k is a random sketching matrix sampled i.i.d at each iteration from a fixed distribution.

Randomized Kaczmarz [13, 28] was the first algorithm of this type. In [11], this sketch-and-project algorithm was analyzed in its full generality. Note that the dual problem of (1.1) takes the form of a quadratic minimization problem [9], and randomized methods such as coordinate descent [18, 31], random pursuit [27, 26] or stochastic dual ascent [9] can thus also be captured as special instances of this method. Richtárik and Takáč [24] adopt a new point of view through a theory of stochastic reformulations of linear systems. In addition, they consider the addition of a relaxation parameter, as well as mini-batch and accelerated variants. Acceleration was only achieved for the expected iterates,

and not in the L2 sense as we do here. We refer to Richtárik and Takáč [24] for interpretation of sketch-and-project as stochastic gradient descent, stochastic Newton descent, stochastic proximal point method, and stochastic fixed point method.

Gower [10] observed that the procedure (1.1) can also be applied to find the inverse of a matrix. Assume the optimization variable itself is a matrix, $x = X$, $b = I$, the identity matrix, then sketch-and-project converges (under mild assumptions) to a solution of $AX = I$. Even the symmetry constraint $X = X^\top$ can be incorporated into the sketch-and-project framework since it is a linear constraint.

There has been recent developments in speeding up the sketch-and-project method using the idea of Nesterov’s acceleration [19]. In [15] an accelerated Kaczmarz algorithm was presented for special sketches of rank one. Arbitrary sketches of rank one were considered in [27], block sketches in [20] and recently, Tu and coauthors [29] developed acceleration for special sketching matrices, assuming the matrix A is square. This assumption, along with any assumptions on A , was later dropped in [23]. Another notable way to accelerate the sketch-and-project algorithm is by using momentum or stochastic momentum [16].

We build on recent work of Richtárik and Takáč [23] and further extend their analysis by studying accelerated sketch-and-project in general Euclidean spaces. This allows us to deduce the result for matrix inversion as a special case. However, there is one additional caveat that has to be considered for the intended application in quasi-Newton methods: ideally, all iterates of the algorithm should be symmetric positive definite matrices. This is not the case in general, but we address this problem by constructing special sketch operators that preserve symmetry and positive definiteness.

1.2 Contributions

We now present our main contributions.

Accelerated Sketch and Project in Euclidean Spaces. We generalize the analysis of an accelerated version of the sketch-and-project algorithm [23] to linear operator systems in Euclidean spaces. We provide a self-contained convergence analysis, recovering the original results in a more general setting.

Faster Algorithms for Matrix Inversion. We develop an accelerated algorithm for inverting positive definite matrices. This algorithm can be seen as a special case of the accelerated sketch-and-project in Euclidean space, thus its convergence follows from the main theorem. However, we also provide a different formulation of the proof that is specialized to this setting.

Similar as in [29], the performance of the algorithm depends on two parameters μ and ν that capture spectral properties of the input matrix and the sketches that are used. Whilst for the non-accelerated sketch-and-project algorithm for matrix inversion [10] the knowledge of these parameters is not necessary, they need to be given as input to the accelerated scheme. When employed with the correct choice of parameters, the accelerated algorithm is always faster than the non-accelerated one. We also provide a theoretical rate for sub-optimal parameters μ, ν , and we perform numerical experiments to argue the choice of μ, ν in practice.

Randomized Accelerated Quasi-Newton. The proposed iterative algorithm for matrix inversion is designed in such a way that each iterate is a symmetric matrix. This means, we can use the generated approximate solutions as estimators for the inverse Hessian in quasi-Newton methods, which is a direct

extension of stochastic quasi-Newton methods. To the best of our knowledge, this yields the first accelerated (stochastic) quasi-Newton method.

Accelerated Quasi-Newton. In the standard BFGS method the updates to the Hessian estimate are not chosen randomly, but deterministically. Based on the intuition gained from the accelerated random method, we propose an accelerated scheme for BFGS. The main idea is that we replace the random sketching of the Hessian with a deterministic update. The theoretical convergence rates do not transfer to this scheme, but we demonstrate by numerical experiments that it is possible to choose a parameter combination which yields a slightly faster convergence. We believe that the novel idea of accelerating BFGS update is extremely valuable, as until now, acceleration techniques were only considered to improve gradient updates.

1.3 Outline

Our accelerated sketch-and-project algorithm for solving linear systems in Euclidean spaces is developed and analyzed in Section 2, and is used later in Section 3 to analyze an accelerated sketch-and-project algorithm for matrix inversion. The accelerated sketch-and-project algorithm for matrix inversion is then used to accelerate the BFGS update, which in turn leads to the development of an accelerated BFGS optimization method. Lastly in Section 4, we perform numerical experiments to gain different insights into the newly developed methods. Proofs of all results and additional insights can be found in the appendix.

2 Accelerated Stochastic Algorithm for Matrix Inversion

In this section we propose an accelerated randomized algorithm to solve linear systems in Euclidean spaces. It will be used later on in order to analyze our newly proposed matrix inversion algorithm, which we then use to estimate the inverse of the Hessian within a quasi-Newton method.

Let \mathcal{X} and \mathcal{Y} be finite dimensional Euclidean spaces and let $\mathcal{A} : \mathcal{X} \mapsto \mathcal{Y}$ be a linear operator. Let $L(\mathcal{X}, \mathcal{Y})$ denote the space of linear operators that map from \mathcal{X} to \mathcal{Y} . Consider the linear system

$$\mathcal{A}x = b, \tag{4}$$

where $x \in \mathcal{X}$ and $b \in \mathbf{Range}(\mathcal{A})$. Consequently there exists a solution to the equation (4). In particular, we aim to find the solution closest to a given initial point $x_0 \in \mathcal{X}$:

$$x^* \stackrel{\text{def}}{=} \arg \min_{x \in \mathcal{X}} \frac{1}{2} \|x - x_0\|^2, \quad \text{subject to } \mathcal{A}x = b. \tag{5}$$

Using the pseudoinverse and Lemma 22 item vi, the solution to (5) is given by

$$x^* = x_0 - \mathcal{A}^\dagger(\mathcal{A}x_0 - b) \in x_0 + \mathbf{Range}(\mathcal{A}^*), \tag{6}$$

where \mathcal{A}^\dagger and \mathcal{A}^* denote the pseudoinverse and the adjoint of \mathcal{A} , respectively.

2.1 The algorithm

Let \mathcal{Z} be a Euclidean space and consider a random linear operator $S_k \in L(\mathcal{Y}, \mathcal{Z})$ chosen from some distribution \mathcal{D} over $L(\mathcal{Y}, \mathcal{Z})$ at iteration k . Our method is given in Algorithm 1, where $Z_k \in L(\mathcal{X})$ is a random linear operator given by the following compositions

$$Z_k = Z(S_k) \stackrel{\text{def}}{=} \mathcal{A}^* S_k^* (S_k \mathcal{A} \mathcal{A}^* S_k^*)^\dagger S_k \mathcal{A}. \quad (7)$$

The updates of variables g_k and x_{k+1} on lines 8 and 9, respectively, correspond to what is known as the *sketch-and-project* update:

$$\begin{aligned} x_{k+1} &= \arg \min_{x \in \mathcal{X}} \frac{1}{2} \|x - y_k\|^2 \\ \text{subject to } & S_k \mathcal{A} x = S_k b, \end{aligned} \quad (8)$$

which can also be written as the following operation

$$x_{k+1} - x_* = (I - Z_k)(y_k - x_*). \quad (9)$$

This can be seen from the fact that $b \in \text{Range}(\mathcal{A})$ together with item i of Lemma 22. Furthermore, note that the adjoint \mathcal{A}^* and the pseudoinverse in Algorithm 1 are taken with respect to the norm in (5).

Algorithm 1 Accelerated Sketch-and-Project (Richtárik and Takáč [23])

- 1: **Parameters:** $\mu, \nu > 0$, \mathcal{D} = distribution over random linear operators.
 - 2: Choose $x_0 \in \mathcal{X}$.
 - 3: Set $v_0 = x_0$
 - 4: Set $\beta = 1 - \sqrt{\frac{\mu}{\nu}}$, $\gamma = \sqrt{\frac{1}{\mu\nu}}$, $\alpha = \frac{1}{1+\gamma\nu}$.
 - 5: **for** $k = 0, 1, \dots$ **do**
 - 6: $y_k = \alpha v_k + (1 - \alpha)x_k$
 - 7: Sample an independent copy $S_k \sim \mathcal{D}$
 - 8: $g_k = \mathcal{A}^* S_k^* (S_k \mathcal{A} \mathcal{A}^* S_k^*)^\dagger S_k (\mathcal{A} y_k - b) = Z_k(y_k - x_*)$
 - 9: $x_{k+1} = y_k - g_k$
 - 10: $v_{k+1} = \beta v_k + (1 - \beta)y_k - \gamma g_k$
 - 11: **end for**
-

Algorithm 1 was first proposed and analyzed by Richtárik and Takáč [23] in the special case when $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{Y} = \mathbb{R}^m$. Our contribution here is in extending the algorithm and analysis to the more abstract setting of Euclidean spaces. In addition, we provide some further extensions of this method in Sections 2.5 and 2.6.

2.2 Key assumptions and quantities

Denote $Z = Z(S)$ for $S \sim \mathcal{D}$. Assume that the *exactness property* holds

$$\begin{aligned} \text{Null}(\mathcal{A}) &= \text{Null}(\mathbf{E}[Z]) \\ \text{Range}(\mathcal{A}^*) &= \text{Range}(\mathbf{E}[Z]). \end{aligned} \quad (10)$$

The exactness assumption is very common in sketch-and-project framework, and indeed it is not very strong. For example, it holds for the matrix inversion problem with every sketching strategies we consider. We further assume that $\mathcal{A} \neq 0$ and $\mathbf{E}[Z]$ is finite. First we collect a few observation on the Z operator

Lemma 1. *The Z operator (7) is a self-adjoint positive projection. Consequently $\mathbf{E}[Z]$ is a self-adjoint positive operator.*

The two parameters that govern the acceleration are

$$\mu \stackrel{\text{def}}{=} \inf_{x \in \text{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z] x, x \rangle}{\langle x, x \rangle} \quad (11)$$

$$\nu \stackrel{\text{def}}{=} \sup_{x \in \text{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z \mathbf{E}[Z]^\dagger Z] x, x \rangle}{\langle \mathbf{E}[Z] x, x \rangle}. \quad (12)$$

The supremum in the definition of ν is well defined due to the exactness assumption together with $\mathcal{A} \neq 0$.

Lemma 2. *We have that*

$$1 \leq \nu \leq \frac{1}{\mu} = \|\mathbf{E}[Z]^\dagger\|. \quad (13)$$

Moreover, if $\text{Range}(\mathcal{A}^*) = \mathcal{X}$, we have

$$\frac{\text{Rank}(\mathcal{A}^*)}{\mathbf{E}[\text{Rank}(Z)]} \leq \nu. \quad (14)$$

2.3 Convergence and change of the norm

For a positive self-adjoint $G \in L(\mathcal{X})$ and $x \in \mathcal{X}$ let $\|x\|_G \stackrel{\text{def}}{=} \sqrt{\langle x, x \rangle_G} \stackrel{\text{def}}{=} \sqrt{\langle Gx, x \rangle}$. We now informally state the convergence rate of Algorithm 1. Theorem 3 generalizes the main theorem from [23] to linear systems in Euclidean spaces.

Theorem 3. *Let x_k, v_k be the random iterates of Algorithm 1. Then*

$$\mathbf{E} \left[\|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \frac{1}{\mu} \|x_k - x_*\|^2 \right] \leq \left(1 - \sqrt{\frac{\mu}{\nu}} \right)^k \mathbf{E} \left[\|v_0 - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \frac{1}{\mu} \|x_0 - x_*\|^2 \right].$$

This theorem shows the accelerated Sketch-and-Project algorithm converges linearly with a rate of $1 - \sqrt{\frac{\mu}{\nu}}$, which translates to a total of $O\left(\sqrt{\frac{\nu}{\mu}} \log\left(\frac{1}{\epsilon}\right)\right)$ iterations to bring the given error in Theorem 3 below $\epsilon > 0$. This is in contrast with the non-accelerated Sketch-and-Project algorithm which requires $O\left(\frac{1}{\mu} \log\left(\frac{1}{\epsilon}\right)\right)$ iterations, as shown in [11] for solving linear systems. From (13) we have that

$$\frac{1}{\sqrt{\mu}} \leq \sqrt{\frac{\nu}{\mu}} \leq \frac{1}{\mu}.$$

On one extreme the above inequality shows that the iteration complexity of the accelerated algorithm is at least as good as its non-accelerated counterpart. On the other extreme, the accelerated algorithm might require as little as the square root of the number of iterations of its non-accelerated counterpart. Since the cost of a single iteration of the accelerated algorithm is of the same order as the non-accelerated algorithm, this theorem shows that acceleration can offer a significant speed-up, which is verified numerically in Section 4.

Note that it is also possible to get the convergence rate of accelerated sketch-and-project where projections are taken with respect to a different weighted norm. For technical details, see Section A.4 of the Appendix.

2.4 Coordinate sketches with convenient probabilities

Let us consider a simple example in the setting for Algorithm 1 where we can understand parameters μ, ν . In particular, consider a linear system $Ax = b$ in \mathbb{R}^n where A is symmetric positive definite.

Corollary 4. Choose $B = A$ and $S = e_i$ with probability proportional to $A_{i,i}$. Then

$$\mu = \frac{\lambda_{\min}(A)}{\text{Tr}(A)} =: \mu^P \quad \text{and} \quad \nu = \frac{\text{Tr}(A)}{\min_i A_{i,i}} =: \nu^P \quad (15)$$

and therefore the convergence rate given in Theorem 3 for the accelerated algorithm is

$$\left(1 - \sqrt{\frac{\mu}{\nu}}\right)^k = \left(1 - \frac{\sqrt{\lambda_{\min}(A) \min_i A_{i,i}}}{\text{Tr}(A)}\right)^k \quad (16)$$

Proof. Clearly, $Z = \frac{1}{A_{i,i}} A^{\frac{1}{2}} S S^{\top} A^{\frac{1}{2}}$, and hence $\mathbf{E}[Z] = \frac{A}{\text{Tr}(A)}$ and $\mu^P = \frac{\lambda_{\min}(A)}{\text{Tr}(A)}$. After simple algebraic manipulations we get

$$\mathbf{E} \left[\mathbf{E}[Z]^{-\frac{1}{2}} Z \mathbf{E}[Z]^{-1} Z \mathbf{E}[Z]^{-\frac{1}{2}} \right] = \text{Tr}(A)^2 \mathbf{E} \left[\frac{1}{A_{i,i}^2} S S^{\top} S S^{\top} \right] = \text{Tr}(A) \text{Diag} \left(A_{i,i}^{-1} \right),$$

$$\text{and therefore } \nu^P = \lambda_{\max} \mathbf{E} \left[\mathbf{E}[Z]^{-\frac{1}{2}} Z \mathbf{E}[Z]^{-1} Z \mathbf{E}[Z]^{-\frac{1}{2}} \right] = \frac{\text{Tr}(A)}{\min_i A_{i,i}}. \quad \square$$

Rate (16) of our accelerated method is to be contrasted to the rate of the non-accelerated method :

$$(1 - \mu)^k = \left(1 - \frac{\lambda_{\min}(A)}{\text{Tr}(A)}\right)^k.$$

Therefore, we gain from acceleration if the smallest diagonal element of A is significantly greater than the smallest eigenvalue, which is a particular type of ill-conditioning.

In fact, parameters μ^P, ν^P above are the correct choice for the matrix inversion algorithm, when symmetry is not enforced, as we shall see later. Unfortunately, we are not able to estimate the parameters while enforcing symmetry for different sketching strategies. We dedicate a section in numerical experiments to test, if the parameter selection (15) performs well under enforced symmetry and different sketching strategies, and also how one might safely choose μ, ν in practice.

2.5 Adding a stepsize ω

In the rest of this section we enrich Algorithm 1 with several *additional* parameters and study their effect on convergence of the resulting method.

First, we consider an extension of Algorithm 1 to a variant which uses a *stepsize parameter* $0 < \omega < 2$. That is, instead of performing the update

$$x_{k+1} = y_k - g_k, \quad (17)$$

we perform the update

$$x_{k+1} = y_k - \omega g_k. \quad (18)$$

Parameters α, β, γ are adjusted accordingly. The resulting method enjoys the rate $\mathcal{O} \left(\left(1 - \sqrt{\frac{\nu}{\mu} \omega (2 - \omega)}\right)^k \right)$, recovering the rate from Theorem 3 as a special case for $\omega = 1$. The formal statement follows.

Theorem 5. Let $0 < \omega < 2$ be an arbitrary stepsize and define

$$\eta \stackrel{\text{def}}{=} 2\omega - \omega^2 \geq 0. \quad (19)$$

Consider a modification of Algorithm 1 where instead of (17) we perform the update (18). If we use the parameters

$$\alpha = \frac{1}{1 + \gamma\nu} \quad \beta = 1 - \sqrt{\frac{\mu\eta}{\nu}} \quad \gamma = \sqrt{\frac{\eta}{\mu\nu}}, \quad (20)$$

then the iterates $\{v_k, x_k\}_{k \geq 0}$ of Algorithm 1 satisfy

$$\mathbf{E} \left[\|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \frac{1}{\mu} \|x_k - x_*\|^2 \right] \leq \left(1 - \sqrt{\frac{\mu\eta}{\nu}} \right)^k \mathbf{E} \left[\|v_0 - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \frac{1}{\mu} \|x_0 - x_*\|^2 \right].$$

Proof. See Appendix B. □

2.6 Allowing for different α

In this section we study how the choice of the key parameter α affects the convergence rate.

This parameter determines how much the sequence $y_k = \alpha v_k + (1 - \alpha)x_k$ resembles the sequence given by x_k or by v_k . For instance, when $\alpha = 0$, $y_k \equiv x_k$, i.e., we recover the steps of the non-accelerated method, and thus one would expect to obtain the same convergence rate as the non-accelerated method. Similar considerations hold in the other extreme, when $\alpha \rightarrow 1$. We investigate this hypothesis, and especially discuss how β and γ must be chosen as a function of α to ensure convergence.

The following statement is a generalization of Theorem 3. For simplicity, we assume that the optional stepsize that was introduced in Theorem 5 is set to one again, $\omega \equiv 1$.

Theorem 6. Let $0 < \alpha < 1$ be fixed. Then the iterates $\{v_k, x_k\}_{k \geq 0}$ of Algorithm 1 with parameters

$$\beta(s) = \frac{1 + s - s\sqrt{\frac{\nu + 4\mu s - 2\nu s + \nu s^2}{\nu s^2}}}{2s}, \quad \gamma(s) = \frac{1}{(1 - s\beta(s))\nu}. \quad (21)$$

where $\tau \stackrel{\text{def}}{=} \frac{1-\alpha}{\alpha}$ and $s \stackrel{\text{def}}{=} \frac{\tau}{\beta\gamma}$, satisfy

$$\mathbf{E} \left[\|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \gamma\tau \|x_k - x_*\|^2 \right] \leq \rho^k \mathbf{E} \left[\|v_0 - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \gamma\tau \|x_0 - x_*\|^2 \right].$$

(or put differently):

$$\mathbf{E} \left[\|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + (1 - \alpha)\gamma \|x_k - x_*\|^2 \right] \leq \rho^k \mathbf{E} \left[\|v_0 - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + (1 - \alpha)\gamma \|x_0 - x_*\|^2 \right].$$

where $\rho = \max\{\beta(s), s\beta(s)\} \leq 1$.

Proof. See Appendix C. □

We can now exemplify a few special parameter settings.

Example 7. For $\alpha = 1$, i.e., if $s \rightarrow 0$, we get the rate $\rho = 1 - \frac{\mu}{\nu}$ with $\beta = 1 - \frac{\mu}{\nu}$, $\gamma = \frac{1}{\nu}$.

Example 8. For $\alpha \rightarrow 0$, i.e., in the limit $s \rightarrow \infty$, we get the rate $\rho = 1 - \frac{\mu}{\nu}$.

Example 9. The rate ρ is minimized for $s = 1$, i.e., $\beta = 1 - \sqrt{\frac{\nu}{\mu}}$ and $\gamma = \sqrt{\frac{1}{\mu\nu}}$; recovering Theorem 3.

The best case, in terms of convergence rate for both non-unit stepsize and a variable parameter choice happened to be the default parameter setup. The non-optimal parameter choice was studied in order to have theoretical guarantees for a wider class of parameters, as in practice one might be forced to rely on sub-optimal / inexact parameter choices.

3 Accelerated Stochastic BFGS Update

The update of the inverse Hessian used with quasi-Newton methods, such as the BFGS method, can be seen as sketch-and-project updates to the linear system $AX = I$, while $X = X^\top$ is enforced. In this section, we present an accelerated version of these updates. We provide two different proofs, one using the view from Theorem 3 and the other using vectorization. By mimicking the updates of the accelerated stochastic BFGS method for inverting matrices, we determine a heuristic for accelerating the classic deterministic BFGS update. We then incorporate this heuristic acceleration into the classic BFGS optimization method and show that the resulting algorithm can offer a speed-up in relation to the standard BFGS algorithm.

3.1 Accelerated matrix inversion

Consider the symmetric positive definite matrix $A \in \mathbb{R}^{n \times n}$ and the following projection problem

$$\begin{aligned} A^{-1} &= \arg \min_X \|X\|_{F(A)}^2 \\ \text{subject to } &AX = I, \quad X = X^\top, \end{aligned} \tag{22}$$

where $\|X\|_{F(A)} \stackrel{\text{def}}{=} \text{Tr}(AX^\top AX) = \|A^{1/2}XA^{1/2}\|_F^2$. This projection problem can be cast as an instantiation of the general projection problem (5). Indeed, we need only note that the constraint in (22) is linear and equivalent to

$$\mathcal{A}(X) \stackrel{\text{def}}{=} \begin{pmatrix} AX \\ X - X^\top \end{pmatrix} = \begin{pmatrix} I \\ 0 \end{pmatrix}.$$

The matrix inversion problem can be efficiently solved using sketch-and-project with a symmetric sketch [12]. The symmetric sketch is given by

$$S_k \mathcal{A}(X) = \begin{pmatrix} S_k^\top AX \\ X - X^\top \end{pmatrix},$$

where $S_k \in \mathbb{R}^{n \times \tau}$ is a random matrix drawn from a distribution \mathcal{D} and $\tau \in \mathbb{N}$. The resulting sketch-and-project method is as follows

$$\begin{aligned} X_{k+1} &= \arg \min_X \|X - X_k\|_{F(A)}^2 \\ \text{subject to } &S_k^\top AX = S_k^\top, \quad X = X^\top, \end{aligned} \tag{23}$$

to which the closed form update is

$$\begin{aligned} X_{k+1} &= S_k(S_k^\top AS_k)^{-1}S_k^\top \\ &+ \left(I - S_k(S_k^\top AS_k)^{-1}S_k^\top A\right) X_k \left(I - AS_k(S_k^\top AS_k)^{-1}S_k^\top\right). \end{aligned} \quad (24)$$

By observing that (3.2) is the sketch-and-project algorithm applied to a linear operator equation, we have constructed an accelerated version in Algorithm 2. We can also apply Theorem 3 to prove that Algorithm 2 is indeed accelerated.

Theorem 10. *The iterates of Algorithm 2 are such that*

$$\begin{aligned} &\mathbf{E} \left[\|V_{k+1} - A^{-1}\|_M^2 + \frac{1}{\mu} \|X_{k+1} - A^{-1}\|_{F(A)}^2 \right] \\ &\leq \left(1 - \sqrt{\frac{\mu}{\nu}}\right) \mathbf{E} \left[\|V_k - A^{-1}\|_M^2 + \frac{1}{\mu} \|X_k - A^{-1}\|_{F(A)}^2 \right], \end{aligned} \quad (25)$$

where $\|X\|_M^2 = \text{Tr} \left(A^{1/2} X^\top A^{1/2} \mathbf{E}[Z]^\dagger A^{1/2} X A^{1/2} \right)$. Furthermore,

$$\begin{aligned} \mu &\stackrel{\text{def}}{=} \inf_{X \in \mathbb{R}^{n \times n}} \frac{\langle \mathbf{E}[Z] X, X \rangle}{\langle X, X \rangle} = \lambda_{\min}(\mathbf{E}[Z]) \\ \nu &\stackrel{\text{def}}{=} \sup_{X \in \mathbb{R}^{n \times n}} \frac{\langle \mathbf{E}[Z \mathbf{E}[Z]^\dagger Z] X, X \rangle}{\langle \mathbf{E}[Z] X, X \rangle}, \end{aligned} \quad (26)$$

where

$$\begin{aligned} \mathbf{Z} &\stackrel{\text{def}}{=} I \otimes I - (I - P) \otimes (I - P), \\ P &\stackrel{\text{def}}{=} A^{1/2} S (S^\top AS)^{-1} S^\top A^{1/2}, \end{aligned} \quad (27)$$

and $Z : X \in \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ is given by $Z(X) = X - (I - P) X (I - P) = XP + PX(I - P)$. Moreover, $2\lambda_{\min}(\mathbf{E}[P]) \geq \lambda_{\min}(\mathbf{E}[\mathbf{Z}]) \geq \lambda_{\min}(\mathbf{E}[P])$.

Notice that preserving the symmetry yields $\mu = \lambda_{\min}(\mathbf{E}[\mathbf{Z}])$ which can be up to twice as large as $\lambda_{\min}(\mathbf{E}[P])$, which is the value of the μ parameter of the method without preserving symmetry. This improved rate when using symmetry is new, and was not present in the algorithm's debut publication [12]. Rather, it was only shown that enforcing symmetry has no effect on the rate.

In terms of parameter estimation, once symmetry is not preserved, we fall back onto the setting from Section 2.4. Unfortunately, we were not able to quantify the effect of enforcing symmetry on the parameter ν .

3.2 Vectorizing – a different insight

Define $\text{Vec} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n^2}$ to be a vectorization operator of column-wise stacking and denote $x \stackrel{\text{def}}{=} \text{Vec}(X)$. It can be shown that the sketch-and-project operation for matrix inversion (3.2) is equivalent to

$$\begin{aligned} x_{k+1} &= \arg \min_x \|x - x_k\|_{A \otimes A}^2 \\ \text{subject to } & (I \otimes S_k^\top)(I \otimes A)x = (I \otimes S_k^\top) \text{Vec}(I), \\ & Cx = 0, \end{aligned}$$

Algorithm 2 Accelerated BFGS inversion

```

1: Parameters:  $\mu, \nu > 0$ ,  $\mathcal{D}$  = distribution over random linear operators.
2: Choose  $X_0 \in \mathcal{X}$ .
3: Set  $V_0 = X_0$ 
4: Set  $\beta = 1 - \sqrt{\frac{\mu}{\nu}}$ ,  $\gamma = \sqrt{\frac{1}{\mu\nu}}$ ,  $\alpha = \frac{1}{1+\gamma\nu}$ .
5: for  $k = 0, 1, \dots$  do
6:    $Y_k = \alpha V_k + (1 - \alpha)X_k$ 
7:   Sample an independent copy  $S \sim \mathcal{D}$ 
8:    $X_{k+1} = Y_k + (Y_k A - I)S(S^\top A S)^{-1}S^\top - S(S^\top A S)^{-1}S^\top A Y_k$ 
9:      $+ S(S^\top A S)^{-1}S^\top A Y_k A S(S^\top A S)^{-1}S^\top$ 
10:   $V_{k+1} = \beta V_k + (1 - \beta)Y_k - \gamma(Y_k - X_{k+1})$ 
11: end for

```

where C is defined so that $Cx = 0$ if and only if $X = X^\top$. The above is a sketch-and-project update for a linear system in \mathbb{R}^{n^2} , which allows to obtain an alternative proof of Theorem 10, without using our results from Euclidean spaces. The details are provided in the appendix.

3.3 Accelerated BFGS as an optimization algorithm

As a tweak in the stochastic BFGS allows for a faster estimation of Hessian inverse and therefore more accurate steps of the method, one might wonder if a equivalent tweak might speed up the standard, deterministic BFGS algorithm. The mentioned tweaked version of standard BFGS is proposed as Algorithm 3. We do not state a convergence theorem for this algorithm or propose to use it as a default solver, but we rather introduce it as a novel idea for accelerating optimization algorithms. We leave theoretical analysis for the future work. For now, we perform several numerical experiments, in order to understand the potential and limitations of this new method.

Algorithm 3 Accelerated BFGS

```

1: Parameters:  $\mu, \nu > 0$ ,  $\mathcal{D}$  = distribution over random linear operators, stepsize  $\eta$ .
2: Choose inversion estimator  $X_0 \in \mathcal{X}$  and starting point  $w_0$ 
3: Set  $V_0 = X_0$ 
4: Set  $\beta = 1 - \sqrt{\frac{\mu}{\nu}}$ ,  $\gamma = \sqrt{\frac{1}{\mu\nu}}$ ,  $\alpha = \frac{1}{1+\gamma\nu}$ .
5: for  $k = 0, 1, \dots$  do
6:    $w_{k+1} = w_k - \eta X_k \nabla f(w_k)$ 
7:    $s_k = w_{k+1} - w_k$ ,  $\zeta_k = \nabla f(w_{k+1}) - \nabla f(w_k)$ 
8:    $Y_k = \alpha V_k + (1 - \alpha)X_k$ 
9:    $X_{k+1} = \frac{\delta_k \delta_k^\top}{\delta_k^\top \zeta_k} + \left(I - \frac{\delta_k \zeta_k^\top}{\delta_k^\top \zeta_k}\right) Y_k \left(I - \frac{\zeta_k \delta_k^\top}{\delta_k^\top \zeta_k}\right)$ 
10:   $V_{k+1} = \beta V_k + (1 - \beta)Y_k - \gamma(Y_k - X_{k+1})$ 
11: end for

```

To better understand Algorithm 3, recall that the BFGS updates an estimate of the inverse Hessian as follows

$$\begin{aligned}
X_{k+1} &= \operatorname{argmin}_X \|X - X_k\|_{F(A)}^2 \\
&\text{subject to } X\delta_k = \zeta_k, X = X^\top,
\end{aligned}$$

where $\delta_k = w_{k+1} - w_k$ and $\zeta_k = \nabla f(w_{k+1}) - \nabla f(w_k)$. The above has the following closed form solution

$$X_{k+1} = \frac{\delta_k \delta_k^\top}{\delta_k^\top \zeta_k} + \left(I - \frac{\delta_k \zeta_k^\top}{\delta_k^\top \zeta_k} \right) X_k \left(I - \frac{\zeta_k \delta_k^\top}{\delta_k^\top \zeta_k} \right).$$

This update appears on line 9 of Algorithm 3 with the difference being that it is applied to a matrix Y_k .

4 Numerical Experiments

We perform extensive numerical experiments to bring an additional insight to both performance and parameter selection of Algorithms 2 and 3. More numerical experiments can be found in Section 5 of the appendix. In the first part of this section, we perform experiments related to Section 3 – we test the accelerated matrix inversion algorithm. In the second part, we perform experiments related to Section 3.3.

4.1 BFGS Matrix Inversion

We consider a problem of inverting a matrix A . Three different sketching strategies are studied: Coordinate sketches with convenient probabilities ($S = e_i$ with probability proportional to $A_{i,i}$), coordinate sketches with uniform probabilities ($S = e_i$ w. p. $\frac{1}{n}$) and Gaussian sketches ($S \sim \mathcal{N}(0, I)$). As matrices to be inverted, we use both artificially generated matrices with the access to the spectrum and also Hessians of ridge regression problems from LIBSVM.

We have shown earlier that μ, ν can be estimated as per (15) for coordinate sketches with convenient probabilities without enforcing symmetry. We use the mentioned parameters for the other sketching strategies and while enforcing the symmetry as well. As in practice, one might not have an access to the exact parameters μ, ν for given sketching strategy, we test sensitivity of the algorithm to parameter choice and also we test possible practical parameter choice (heuristic) – when ν is chosen by (15) and $\mu = \frac{1}{100\nu}$ or $\mu = \frac{1}{10000\nu}$.

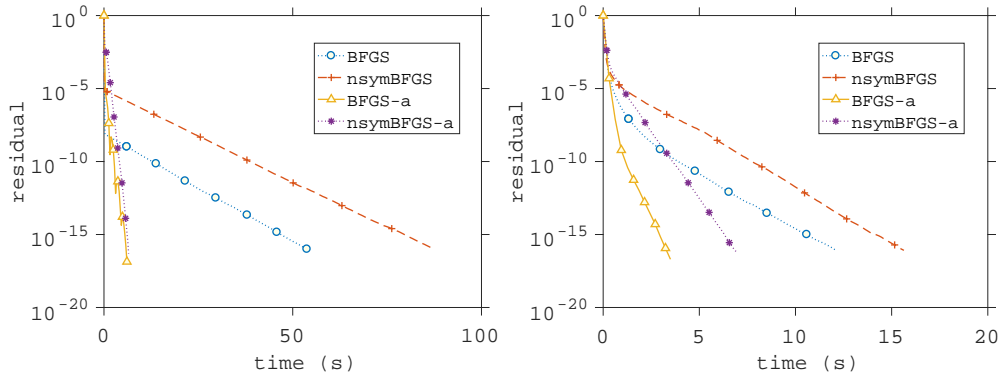


Figure 1: Accelerated algorithm applied on artificial data. Left figure: Eigenvalues of $A \in \mathbb{R}^{100 \times 100}$ are $1, 10^3, 10^3, \dots, 10^3$ and coordinate sketches with convenient probabilities are used. Right figure: Eigenvalues of $A \in \mathbb{R}^{100 \times 100}$ are $1, 2, \dots, n$ and Gaussian sketches are used. Label “nsym” indicates non-enforcing symmetry and “-a” indicates acceleration.

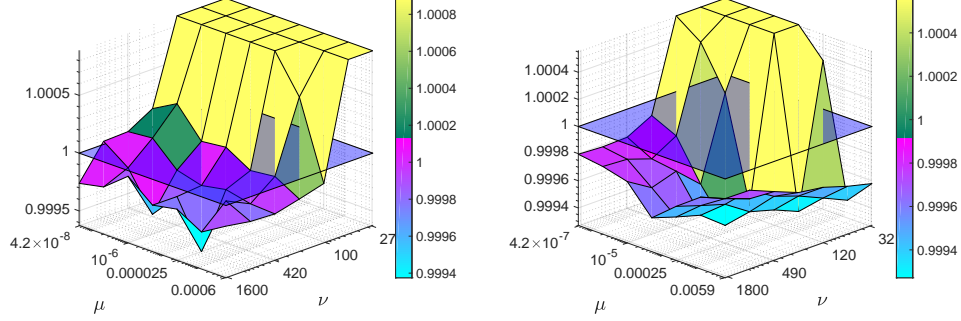


Figure 2: Sensitivity to acceleration parameters. Left figure: Eigenvalues of $A \in \mathbb{R}^{200 \times 200}$ are $1, 10^3, 10^3, \dots, 10^3$ and coordinate sketches with convenient probabilities are used. Right figure: Eigenvalues of $A \in \mathbb{R}^{200 \times 200}$ are $1, 2, \dots, n$ and Gaussian sketches are used. Height corresponds to average per iteration decrease of the residual divided by average per iteration change of residual of nonaccelerated algorithm (the transparent plane). Yellow flat region indicates divergence of the algorithm. Choice of parameters as per (15) in the middle of the plot.

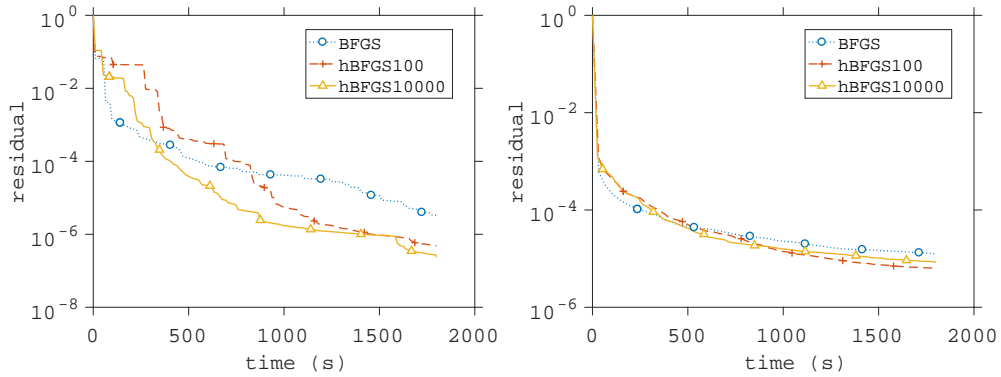


Figure 3: Comparison of nonaccelerated method to accelerated method with heuristic choice of parameters on LIBSVM dataset. Left figure: Epsilon dataset ($n = 2000$), coordinate sketches with uniform probabilities. Right figure: SVHN dataset ($n = 3072$), coordinate sketches with convenient probabilities. Label “h” indicates that λ_{\min} was not precomputed, but μ was chosen as described in the text.

For more plots, see Section 5 in the appendix as here we provide only a tiny fraction of all plots. The experiments suggest that once the parameters μ, ν are estimated exactly, we get a speedup comparing to the nonaccelerated method; and the amount of speedup depends on the structure of A and the sketching strategy. We observe from Figure 1 that we gain a great speedup for ill conditioned problems once the

eigenvalues are concentrated around the largest eigenvalue. We also observe from Figures 1 and 3 that enforcing symmetry combines well with μ, ν computed for the algorithm which do not enforce symmetry. On top of that, choice of μ, ν per (15) seems to be robust to different sketching strategies, and in worst case performs as fast as nonaccelerated algorithm. The sensitivity plot in Figure 2 indicates that the algorithm might even diverge once parameter ν is chosen to be small enough, and is not very sensitive on the choice of μ .

4.2 BFGS Optimization Method

We test Algorithm 3 on several logistic regression problems using data from LIBSVM [4]. In all our tests we centered and normalized the data, included a bias term (a linear intercept), and choose the regularization parameter as $\lambda = 1/m$, where m is the number of data points. To keep things as simple as possible, we also used a fixed stepsize which was determined using grid search. Since our theory regarding the choice for the parameters μ and ν does not apply in this setting, we simply probed the space of parameters manually and reported the best found result, see Figures 4, 5, 6 and 7 where we tested four small problems based on the data sets phishing, mushrooms, australian and splice, respectively. In the legend of these figures we use BFGS-a- μ - ν to denote the accelerated BFGS method (Algorithm 3) with parameters μ and ν .

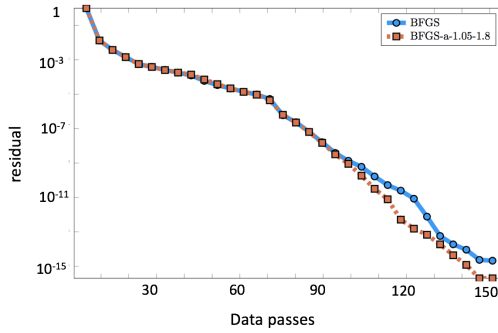


Figure 4: phishing

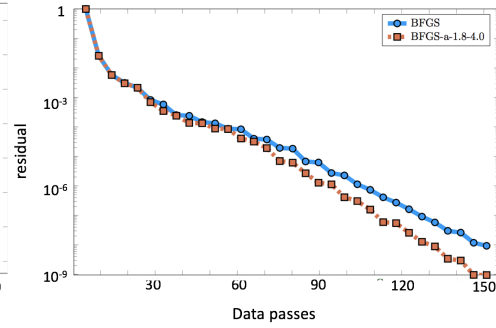


Figure 5: mushrooms

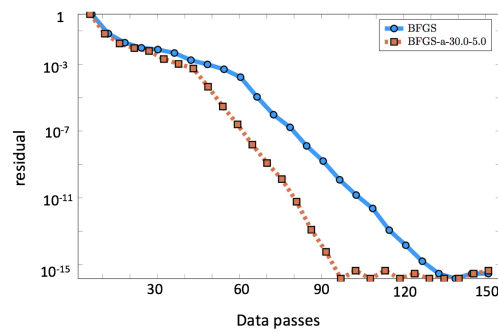


Figure 6: australian

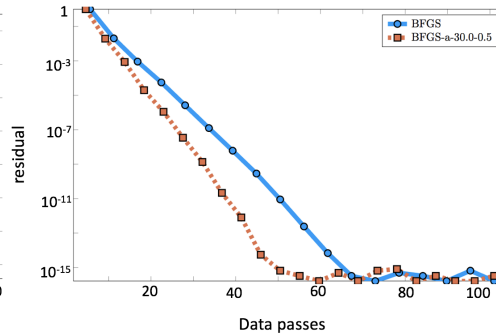


Figure 7: splice

We also give additional experiments with the same setup to the ones found in Section 4.2. In Figure 8 we show a example where acceleration often hinders the performance of the BFGS method. While in

Figures 9 and 10 we show two problems where acceleration has practically no affect. Indeed, we found in our experiments that even when choosing extreme values of μ and ν , the generated inverse Hessian would not significantly deviate from the estimate that one would obtain using the standard BFGS update. Thus on these two problems there is apparently no room for improvement by using acceleration.

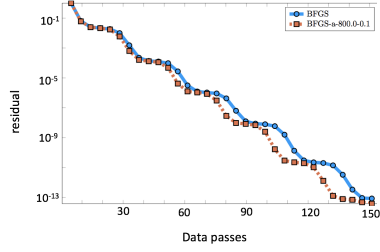


Figure 8: Dataset madelon:

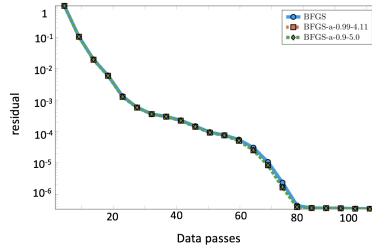


Figure 9: Dataset covtype

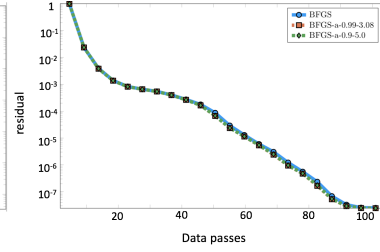


Figure 10: Dataset a9a

5 Further Experiments with Accelerated quasi-Newton Updates

In this section, we test the empirical rate of convergence of Algorithm 2, the accelerated BFGS update for inverting positive definite matrices. Only vector sketches are considered, as the standard quasi-Newton methods update inverse Hessian according to the action in only one direction. We compare speed of accelerated method with precomputed estimate of parameters μ, ν to nonaccelerated method. The precomputed estimate of μ^P, ν^P is set as per (15):

$$\mu^P = \frac{\lambda_{\min}(A)}{\text{Tr}(A)}, \quad \nu^P = \frac{\text{Tr}(A)}{\min_i(A_{i,i})},$$

which is the optimal choice for coordinate sketches with convenient probabilities without enforcing symmetry. As in practice, we might not have an access to $\lambda_{\min}(A)$, thus we cannot compute μ^P exactly. Therefore we also test sensitivity of the algorithm to the choice of parameters, and we run some experiments where we only guess parameter μ^P .

Lastly the tests are performed on both artificial examples and LIBSVM [4] data. We shall also explain the legend of plots – “a” indicates acceleration, “nsym” indicates the algorithm without enforcing symmetry and “h” indicates the setting when ν^P is not known, and a naive heuristic choice is casted.

5.1 Simple and well understood artificial example

Let us consider inverting matrix $A = \alpha I + \beta \mathbf{1}\mathbf{1}^\top$ for $\alpha > 0$ and $\beta \geq -\frac{\alpha}{n}$ so as in this case we have both control over μ and ν . This artificial example was considered in [29] for solving linear systems. In particular, we show that for coordinate sketches with convenient probabilities (which is indeed the same as uniform probabilities in this example), we have

$$\begin{aligned} \mu^P &\stackrel{\text{def}}{=} \lambda_{\min}(\mathbf{E}[P]) = \frac{\min(\alpha, \alpha + n\beta)}{n(\alpha + \beta)}, \\ \nu^P &\stackrel{\text{def}}{=} \lambda_{\max}\left(\mathbf{E}\left[P^{-\frac{1}{2}} P \mathbf{E}[P]^{-1} P P^{-\frac{1}{2}}\right]\right) = n. \end{aligned}$$

Due to the fact that we do not have a theoretical justification of μ, ν for $n > 2$ when enforcing symmetry, we set $\mu = \mu^P$ and $\nu = \nu^P$ for Gaussian sketches as well.

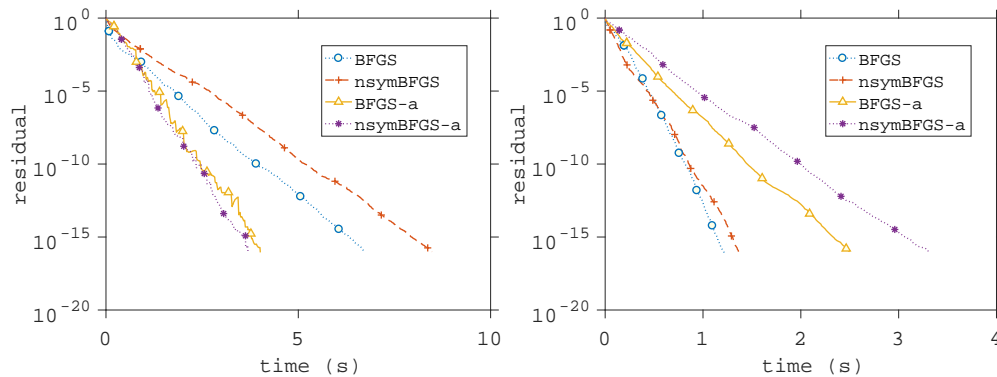


Figure 11: Parameter choice: $\alpha = 1 + 10^{-1}, \beta = -n^{-1}, n = 100$. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch respectively.

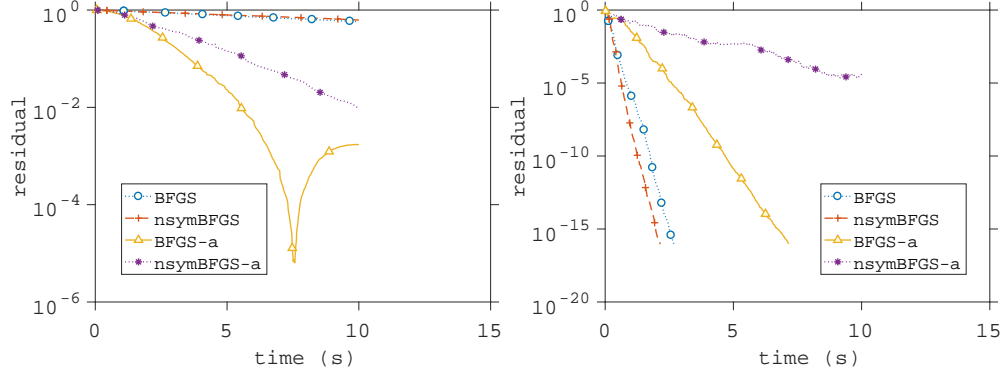


Figure 12: Parameter choice: $\alpha = 1 + 10^{-3}, \beta = -n^{-1}, n = 100$. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch respectively.

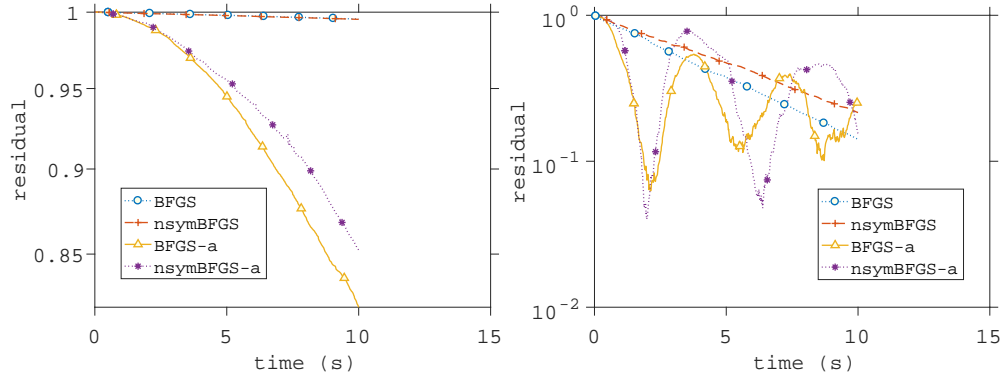


Figure 13: Parameter choice: $\alpha = 1 + 10^{-5}, \beta = -n^{-1}, n = 100$. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch, respectively.

As expected from the theory, as the matrix to be inverted becomes more ill conditioned, the accelerated method performs significantly better comparing to nonaccelerated method for coordinate sketches. In fact, an arbitrary speedup can be obtained by setting $\beta = -n^{-1}$ and $\alpha \rightarrow 1$ for coordinate sketches setup. On the other hand, gaussian sketches report the slowing of the algorithm, most likely caused by the fact that theoretical parameters μ, ν for Gaussian sketches with enforced symmetry are different to μ^P, ν^P , which are estimated for coordinate sketches without enforced symmetry. In the case of coordinate sketches with symmetry enforced, we suspect a great speedup even though the parameters μ, ν were set to μ^P, ν^P .

5.2 Random artificial example

We randomly generate orthonormal matrix U , choose diagonal matrix D , and set $A = UDU^\top$. Clearly, diagonal elements of D are eigenvalues of A . We set them in the following way:

- Uniform grid. The eigenvalues are set to $1, 2, \dots, n$.
- One small, rest larger. The smallest eigenvalue is 1, remaining eigenvalues are all 10 in first example, all 100 in second example and all 1000 in the third example in this category.

- One large, rest small. The largest eigenvalue is 10^4 , remaining eigenvalues are all 1.

Firstly, consider coordinate sketches with convenient probabilities. Notice that we can easily estimate ν^P, μ^P due to results from Section 2.4 since we have control of $\lambda_{\min}(A)$ and therefore about μ . Therefore, we set $\mu = \mu^P = \min D_{i,i}$ and $\nu = \nu^P$ for Algorithm 2. Then, we consider coordinate sketches with uniform probabilities and Gaussian sketches. In both cases, we set parameters μ, ν as for coordinate sketches with convenient probabilities.

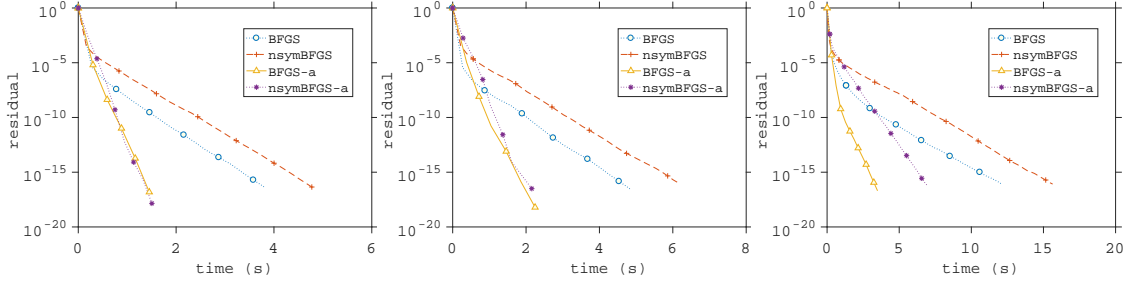


Figure 14: Eigenvalues set to $1, 2, 3, \dots, n$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

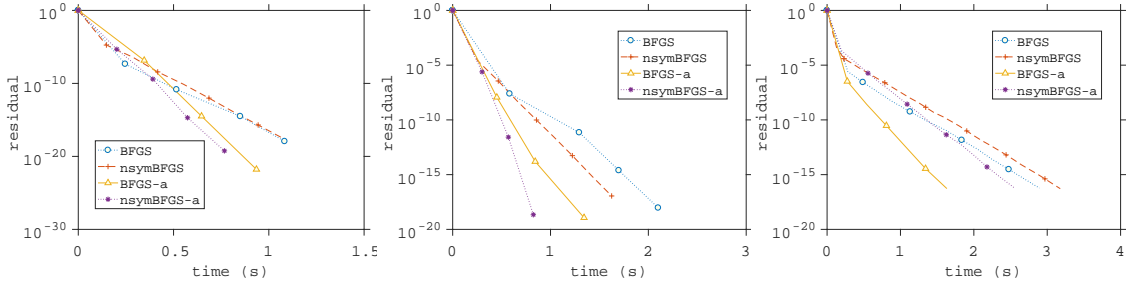


Figure 15: Eigenvalues set to $1, 10, 10, \dots, 10$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

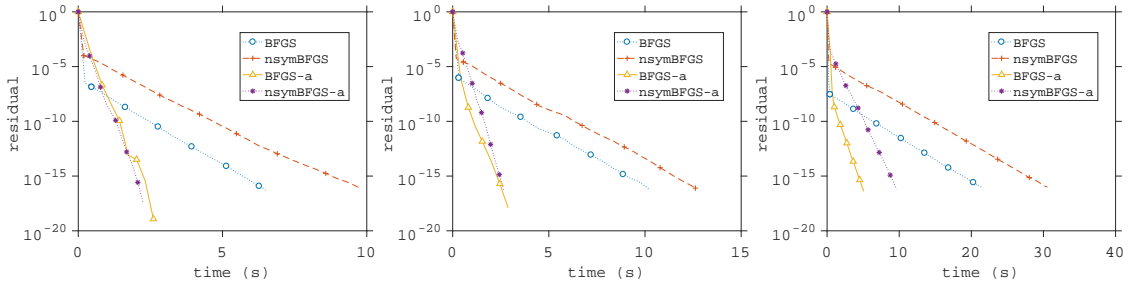


Figure 16: Eigenvalues set to $1, 100, 100, \dots, 100$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

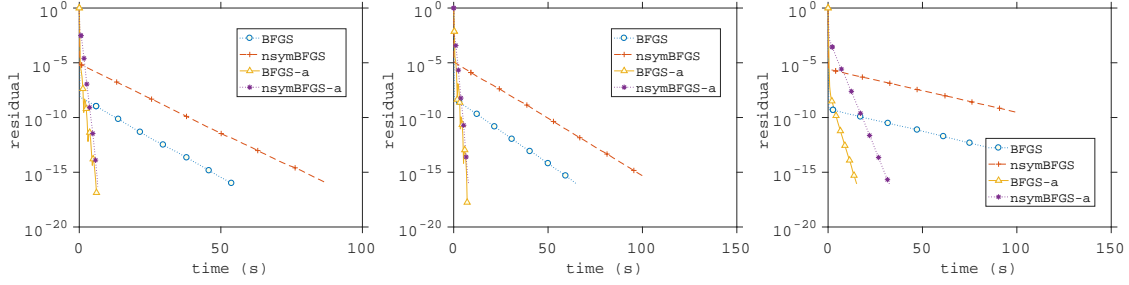


Figure 17: Eigenvalues set to 1, 1000, 1000, \dots 1000. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

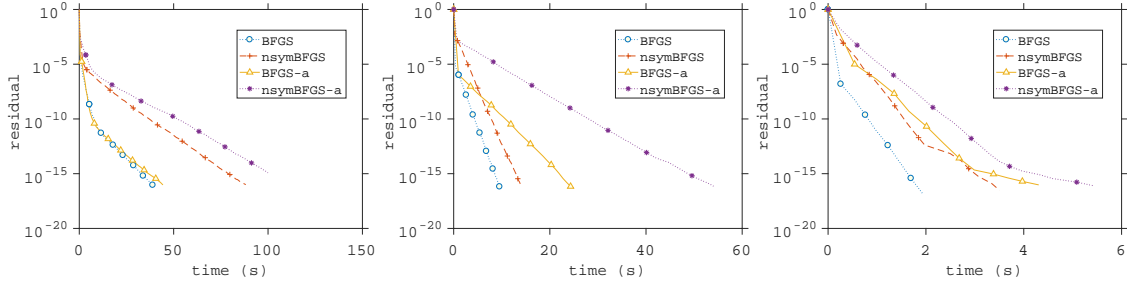


Figure 18: Eigenvalues set to 10000, 1, 1, \dots 1. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

The numerical experiments in this section indicate that one might choose μ, ν as per Section 2.4. In other words, one might pretend to be in the setting when symmetry is not enforced and coordinate sketches with convenient probabilities are used. In fact, the practical speedup coming from the acceleration depends very strongly on the structure of matrix A . Another message to be delivered is that both preserving symmetry and acceleration yield a better convergence and they combine together well.

We also consider a problem where we pretend to not have access to $\lambda_{\min}(A)$, therefore we cannot choose $\mu = \mu^P$. Instead, we naively choose $\mu = \frac{1}{100\nu}$ and $\mu = \frac{1}{10000\nu}$.

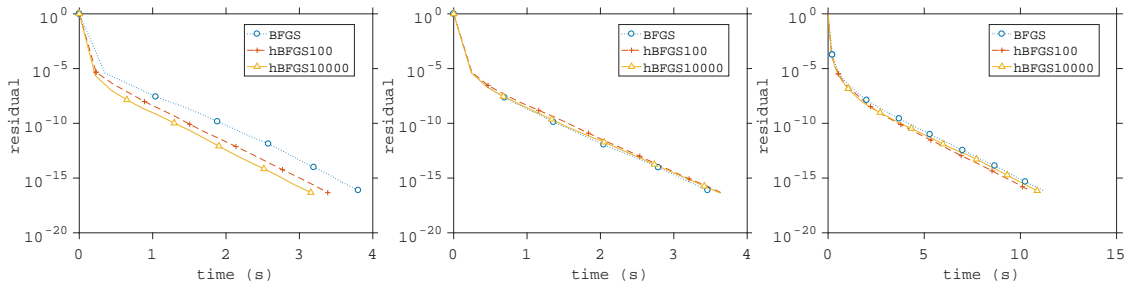


Figure 19: Eigenvalues set to 1, 2, \dots , n . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

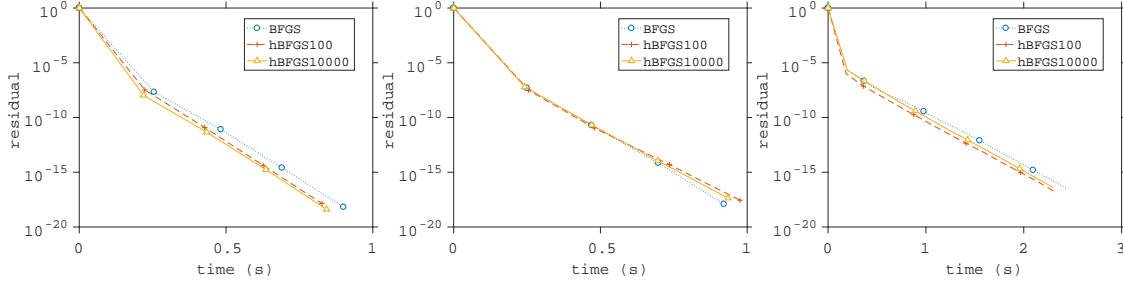


Figure 20: Eigenvalues set to 1, 10, 10, \dots 10. Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

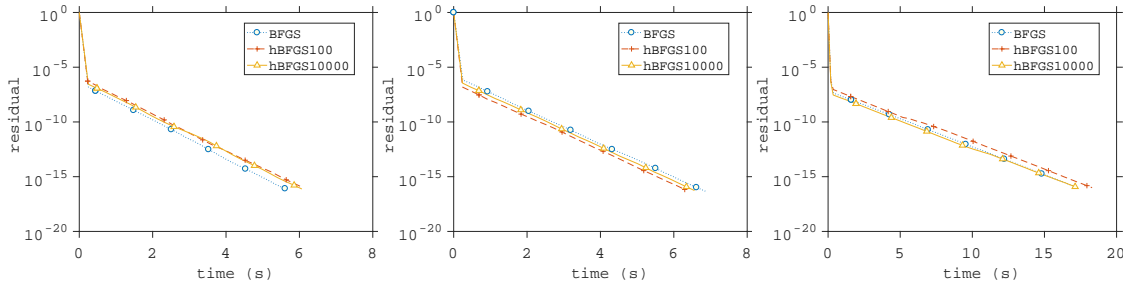


Figure 21: Eigenvalues set to 1, 100, 100, \dots 100. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

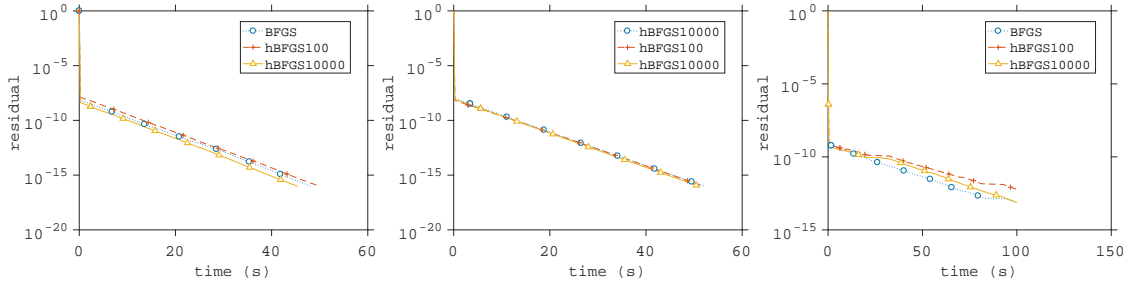


Figure 22: Eigenvalues set to 1, 1000, 1000, \dots 1000. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

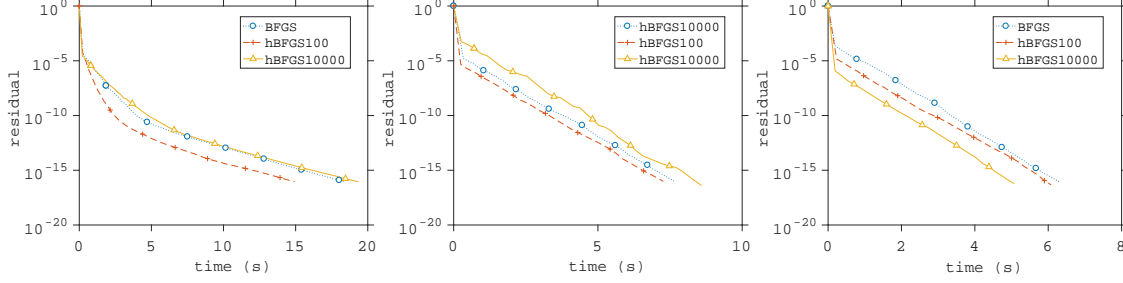


Figure 23: Eigenvalues set to $10000, 1, 1, \dots, 1$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

Notice that once the acceleration parameters are not set exactly (but they are still reasonable), we observe that the performance of the accelerated algorithm is essentially the same as the performance of nonaccelerated algorithm. We have observed the similar behavior when setting $\mu = \mu^P$ for Gaussian sketches.

5.2.1 Sensitivity to the acceleration parameters

Here we investigate the sensitivity of the accelerated BFGS to the parameters μ and ν . First we compute ν^P, μ^P and from this we extract the following exponential grids: $\mu_i = 2^{i-4}\mu$ and $\nu_i = 5^{i-4}\nu$ for $i = 1, 2, \dots, 7$. To gauge the gain is using acceleration with a particular (μ, ν) pair, we run the accelerated algorithm for a fixed time then store the error of the final iterate. We then compute average per iteration decrease and divide it by average per iteration decrease of nonaccelerated algorithm. Thus if the resulting difference is less than one, then the accelerated algorithm was faster to nonaccelerated.

In the plots below, $n = 200$ was chosen. We focused on 2 problems described in the previous section – when the eigenvalues are uniformly distributed and when the the largest eigenvalue have multiplicity $n - 1$.

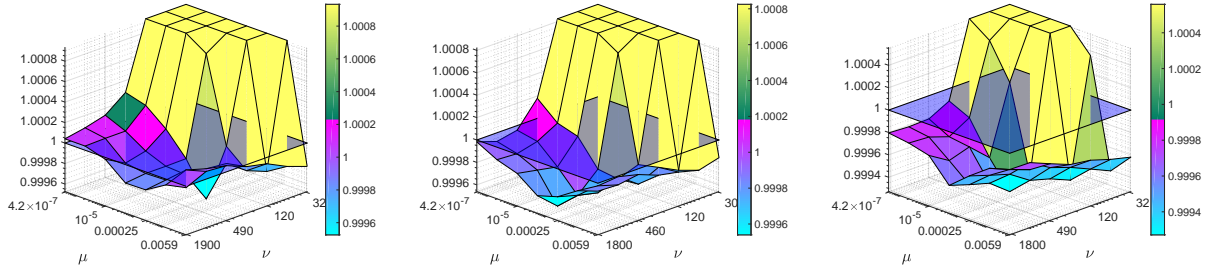


Figure 24: Sensitivity to acceleration parameters. Eigenvalues of A are set to $1, 2, \dots, n$. From left to right we have: Coordinate sketches with convenient probabilities, coordiante sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (15) in the middle of plots. Each instance was run for 5 seconds.

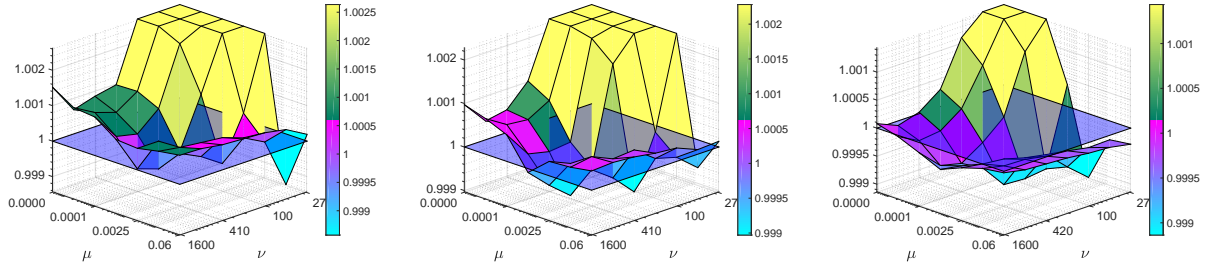


Figure 25: Sensitivity to acceleration parameters. Eigenvalues of A are set to $1, 10, 10, \dots, 10$. From left to right we have: Coordinate sketches with convenient probabilities, coordiante sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (15) in the middle of plots. Each instance was run for 2 seconds.

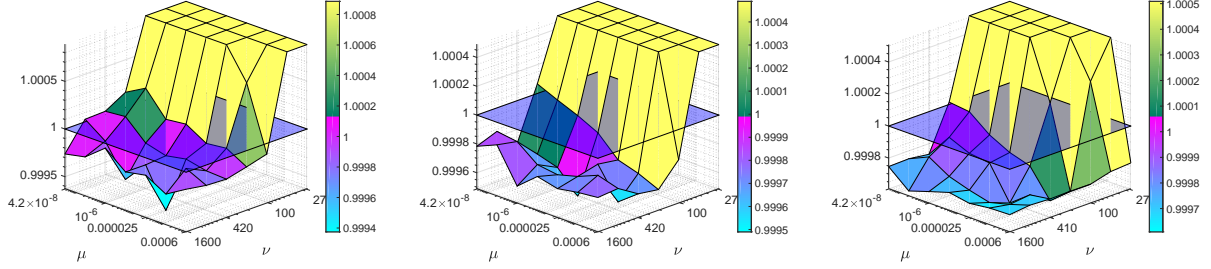


Figure 26: Sensitivity to acceleration parameters. Eigenvalues of A are set to 1, 1000, 1000, \dots , 1000. From left to right we have: Coordinate sketches with convenient probabilities, coordiante sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (15) in the middle of plots. Each instance was run for 10 seconds.

The crucial aspect to make the accelerated algorithm to converge is to set ν large enough. In fact, combination of both small ν and small μ leads almost always to non-convergent algorithm. On the other hand, it seems that once ν is chosen correctly, big enough μ leads to fast convergence. This indicates how to compute μ in practice (recall that computing ν is feasible) – one needs just to choose it small enough (definitely smaller than $\frac{1}{\nu}$).

5.3 Experiments with LIBSVM

Next we investigate if the accelerated BFGS update improves upon the standard BFGS update when applied to the Hessian $\nabla^2 f(x)$ of ridge regression problems of the form

$$\min_{x \in \mathbb{R}^n} f(x) \stackrel{\text{def}}{=} \frac{1}{2} \|Ax - b\|_2^2 + \frac{\lambda}{2} \|x\|_2^2, \quad \nabla^2 f(x) = A^\top A + \lambda I, \quad (28)$$

using data from LIBSVM [4]. Datapoints (rows of A) were normalized such that $\|A_{i,:}\|^2 = 1$ for all i and the regularization parameter was chosen as $\lambda = \frac{1}{m}$.

First, we run the experiments on smaller problems when parameters μ , ν are precomputed for coordinate sketches with convenient probabilities (15).

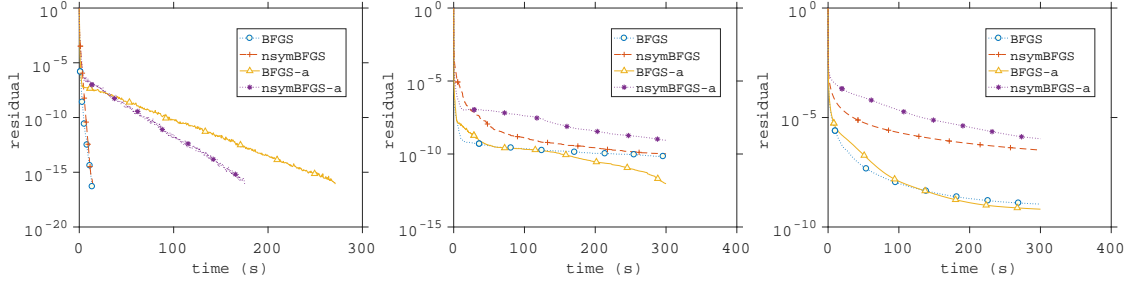


Figure 27: Dataset aloi: $n = 128$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

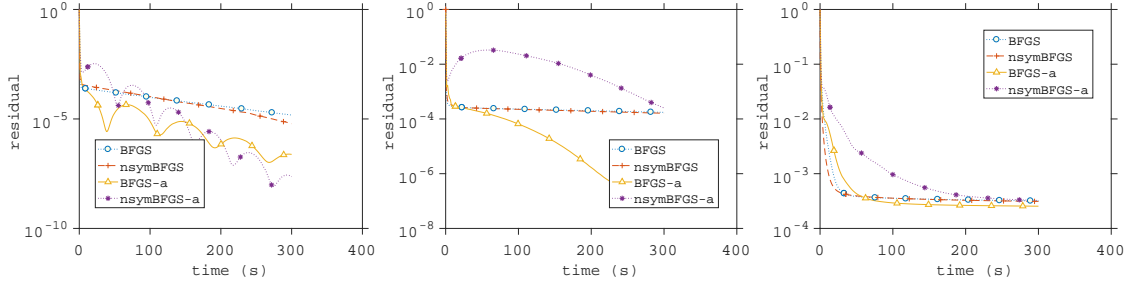


Figure 28: Dataset w1a: $n = 300$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

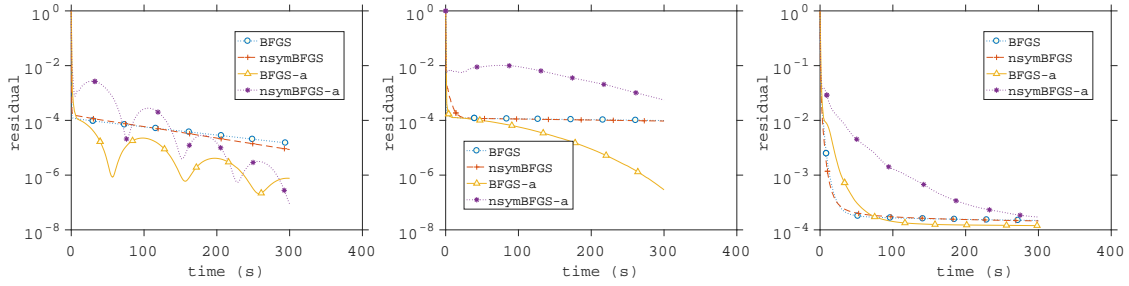


Figure 29: Dataset w2a: $n = 300$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

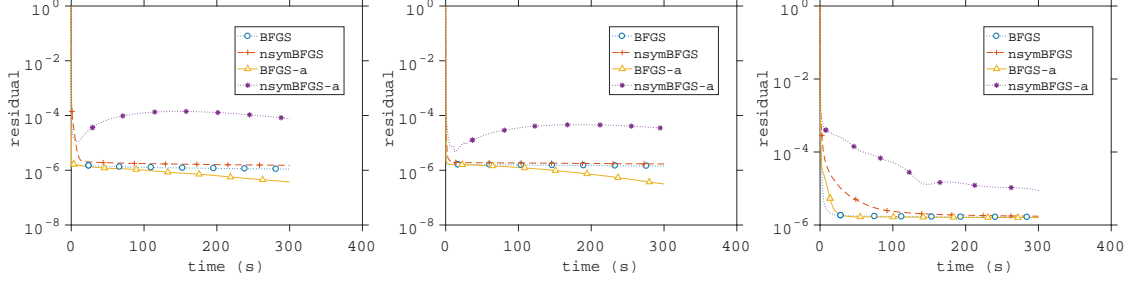


Figure 30: Dataset mushrooms: $n = 112$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

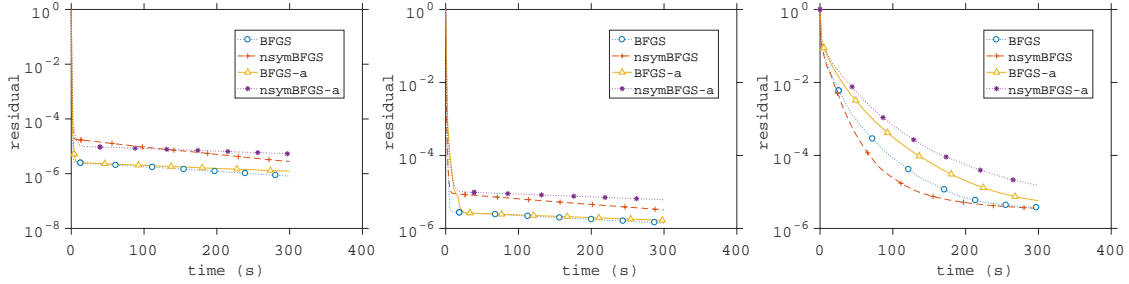


Figure 31: Dataset protein: $n = 357$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

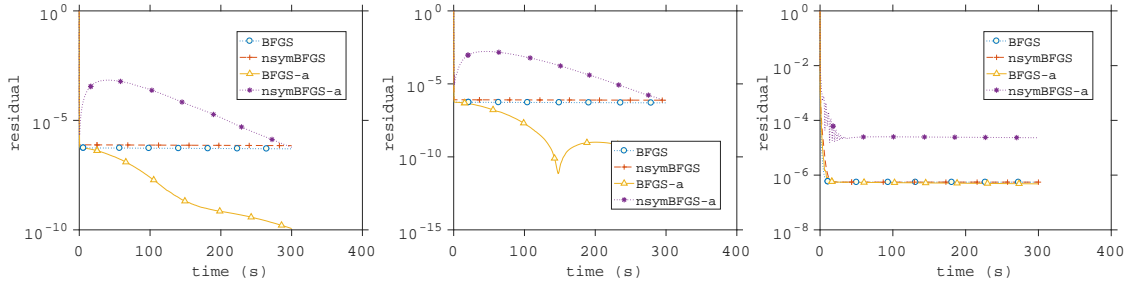


Figure 32: Dataset phishing: $n = 68$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

In vast majority of examples, accelerated method performed significantly better than nonaccelerated method for coordinate sketches (with both convenient and uniform probabilities), however the methods were comparable for Gaussian sketches. We believe that this is due to the fact that choice of parameters as per (15) is close to optimal parameters for coordinate sketches, and further for Gaussian sketches. However, the experiments on coordinate sketches indicates that for some many classes of problems, accelerated algorithm with finely tuned parameters brings a great speedup comparing to nonaccelerated one.

We also consider a problem where we do not compute $\lambda_{\min}(A)$, and therefore we cannot choose $\mu = \mu^P$ in (15). Instead, we choose $\mu = \frac{1}{100\nu}$ and $\mu = \frac{1}{10000\nu}$.

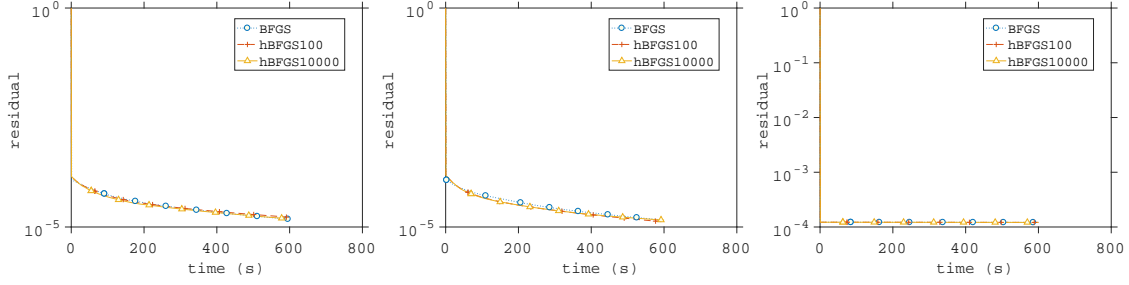


Figure 33: Dataset madelon: $n = 500$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

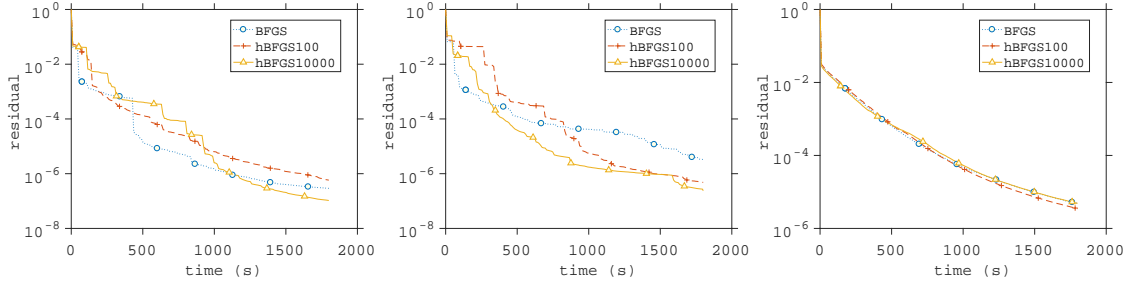


Figure 34: Dataset epsilon: $n = 2000$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

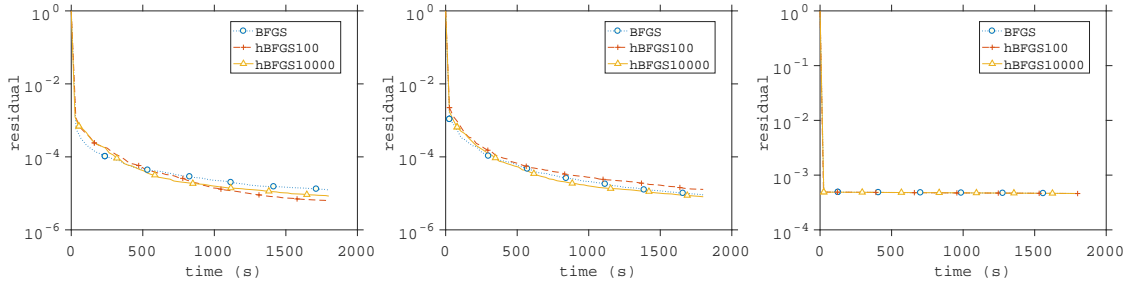


Figure 35: Dataset svhn: $n = 3072$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

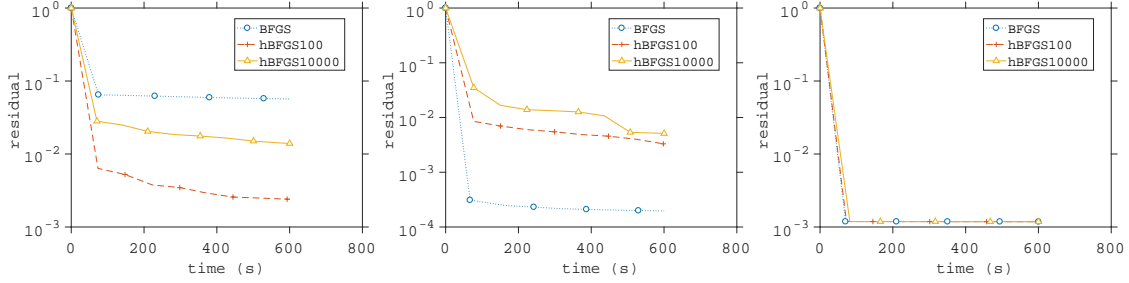


Figure 36: Dataset gisette: $n = 5000$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

Notice that once the acceleration parameters are not set exactly (but they are still reasonable), we observe that the performance of the accelerated algorithm is essentially the same as the performance of nonaccelerated algorithm, which is essentially the same conclusion as for artificially generated examples.

6 Conclusions and Extensions

We developed an accelerated sketch-and-project method for solving linear systems in Euclidean spaces. The method was applied to invert positive definite matrices, while keeping their symmetric structure. Our accelerated matrix inversion algorithm was then incorporated into an optimization framework to develop both accelerated stochastic and deterministic BFGS, which to the best of our knowledge, are *the first accelerated quasi-Newton updates*.

We show that under a careful choice of the parameters of the method, and depending on the problem structure and conditioning, acceleration might result into significant speedups both for the matrix inversion problem and for the stochastic BFGS algorithm. We confirm experimentally that our accelerated methods can lead to speed-ups when compared to the classical BFGS algorithm.

As a future line of research, it might be interesting to study the accelerated BFGS algorithm (either deterministic or stochastic) further, and provide a convergence analysis on a suitable class of functions. Another interesting area of research might be to combine accelerated BFGS with limited memory [14] or engineer the method so that it can efficiently compete with first order algorithms for some empirical risk minimization problems, such as, for example [8].

As we show in this work, *Nesterov's acceleration can be applied to quasi-Newton updates*. We believe this is a surprising fact, as quasi-Newton updates have not been understood as optimization algorithms, which prevented the idea of applying acceleration in this context.

Since second-order methods are becoming more and more ubiquitous in machine learning and data science, we hope that our work will motivate further advances at the frontiers of big data optimization.

References

- [1] Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization for machine learning in linear time. *The Journal of Machine Learning Research*, 18(1):4148–4187, 2017.
- [2] Albert S. Berahas, Raghu Bollapragada, and Jorge Nocedal. An investigation of Newton-sketch and subsampled Newton methods. *CoRR*, abs/1705.06211, 2017.

- [3] Charles G Broyden. Quasi-Newton methods and their application to function minimisation. *Mathematics of Computation*, 21(99):368–381, 1967.
- [4] Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM Trans. Intell. Syst. Technol.*, 2(3):27:1–27:27, May 2011.
- [5] C. A. Desoer and B. H. Whalen. A note on pseudoinverses. *Journal of the Society of Industrial and Applied Mathematics*, 11(2):442–447, 1963.
- [6] Roger Fletcher. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–322, 1970.
- [7] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of computation*, 24(109):23–26, 1970.
- [8] Robert Gower, Donald Goldfarb, and Peter Richtárik. Stochastic block BFGS: Squeezing more curvature out of data. In *International Conference on Machine Learning*, pages 1869–1878, 2016.
- [9] Robert M. Gower and Peter Richtárik. Stochastic dual ascent for solving linear systems. *arXiv:1512.06890*, 2015.
- [10] Robert M. Gower and Peter Richtárik. Randomized quasi-Newton updates are linearly convergent matrix inversion algorithms. *SIAM Journal on Matrix Analysis and Applications*, 38(4):1380–1409, 2017.
- [11] Robert Mansel Gower and Peter Richtárik. Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.
- [12] Robert Mansel Gower and Peter Richtárik. Randomized quasi-Newton updates are linearly convergent matrix inversion algorithms. *SIAM Journal on Matrix Analysis and Applications*, 38(4):pp. 1380–1409, 2017.
- [13] S. Kaczmarz. Angenäherte Auflösung von Systemen linearer Gleichungen. *Bulletin International de l’Académie Polonaise des Sciences et des Lettres*, 35:355–357, 1937.
- [14] Dong C Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1-3):503–528, 1989.
- [15] Ji Liu and Stephen J. Wright. An accelerated randomized Kaczmarz algorithm. *Math. Comput.*, 85(297):153–178, 2016.
- [16] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. *arXiv preprint arXiv:1712.09677*, 2017.
- [17] Philipp Moritz, Robert Nishihara, and Michael Jordan. A linearly-convergent stochastic L-BFGS algorithm. In *Artificial Intelligence and Statistics*, pages 249–258, 2016.
- [18] Yu. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [19] Yurii Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$. *Soviet Mathematics Doklady*, 27(2):372–376, 1983.

- [20] Yurii Nesterov and Sebastian U. Stich. Efficiency of the accelerated coordinate descent method on structured optimization problems. *SIAM Journal on Optimization*, 27(1):110–123, 2017.
- [21] G.K. Pedersen. *Analysis Now*. Graduate Texts in Mathematics. Springer New York, 1996.
- [22] Mert Pilanci and Martin J. Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. *SIAM Journal on Optimization*, 27(1):205–245, 2017.
- [23] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: accelerated method. *Manuscript, October 2017*, 2017.
- [24] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: algorithms and convergence theory. *arXiv:1706.01108*, 2017.
- [25] David F Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathematics of computation*, 24(111):647–656, 1970.
- [26] S. U. Stich, C. L. Müller, and B. Gärtner. Variable metric random pursuit. *Mathematical Programming*, 156(1):549–579, Mar 2016.
- [27] Sebastian U. Stich. *Convex Optimization with Random Pursuit*. PhD thesis, ETH Zurich, 2014. Diss., Eidgenössische Technische Hochschule ETH Zurich, Nr. 22111.
- [28] Thomas Strohmer and Roman Vershynin. A randomized Kaczmarz algorithm with exponential convergence. *Journal of Fourier Analysis and Applications*, 15(2):262, 2009.
- [29] Stephen Tu, Shivaram Venkataraman, Ashia C. Wilson, Alex Gittens, Michael I. Jordan, and Benjamin Recht. Breaking locality accelerates block Gauss-Seidel. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 3482–3491, 2017.
- [30] Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-Newton methods for nonconvex stochastic optimization. *SIAM Journal on Optimization*, 27(2):927–956, 2017.
- [31] Stephen J. Wright. Coordinate descent algorithms. *Math. Program.*, 151(1):3–34, June 2015.
- [32] Peng Xu, Jiyan Yang, Farbod Roosta-Khorasani, Christopher Ré, and Michael W Mahoney. Sub-sampled Newton methods with non-uniform sampling. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3000–3008. Curran Associates, Inc., 2016.

A Proofs for Section 2

A.1 Proof of Lemma 2

First note that Z is a self-adjoint positive operator and thus so is $\mathbf{E}[Z]$. Consequently,

$$\begin{aligned}
 \mu & \stackrel{(11)}{=} \inf_{x \in \text{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z] x, x \rangle}{\langle x, x \rangle} \\
 & \stackrel{(10)}{=} \inf_{x \in \text{Range}(\mathbf{E}[Z])} \frac{\langle \mathbf{E}[Z] x, x \rangle}{\langle x, x \rangle} \\
 & \stackrel{\text{Lemma 22 item ii}}{=} \inf_{x \in \mathcal{X}} \frac{\langle \mathbf{E}[Z] \mathbf{E}[Z]^\dagger x, \mathbf{E}[Z]^\dagger x \rangle}{\langle \mathbf{E}[Z]^\dagger x, \mathbf{E}[Z]^\dagger x \rangle} \\
 & \stackrel{\text{Lemma 22 item i}}{=} \inf_{x \in \mathcal{X}} \frac{\langle \mathbf{E}[Z]^\dagger x, x \rangle}{\langle \mathbf{E}[Z]^\dagger x, \mathbf{E}[Z]^\dagger x \rangle} \\
 & \stackrel{\text{Lemma 18}}{=} \inf_{z \in \text{Range}((\mathbf{E}[Z]^\dagger)^{1/2})} \frac{\langle z, z \rangle}{\langle \mathbf{E}[Z]^\dagger z, z \rangle} \quad (\text{set } z = (\mathbf{E}[Z]^\dagger)^{1/2} x) \\
 & \stackrel{(69)}{=} \frac{1}{\|\mathbf{E}[Z]^\dagger\|}. \tag{29}
 \end{aligned}$$

For the bounds (13) we have that

$$\begin{aligned}
 \nu & \stackrel{(12)}{=} \sup_{x \in \text{Range}(\mathcal{A}^*)} \frac{\mathbf{E} \left[\langle \mathbf{E}[Z]^\dagger Z x, Z x \rangle \right]}{\langle \mathbf{E}[Z] x, x \rangle} \\
 & \leq \sup_{x \in \text{Range}(\mathcal{A}^*)} \frac{\|\mathbf{E}[Z]^\dagger\| \mathbf{E} [\|Z x\|_2^2]}{\langle \mathbf{E}[Z] x, x \rangle} \\
 & = \|\mathbf{E}[Z]^\dagger\| \\
 & \stackrel{(29)}{\leq} \frac{1}{\mu}.
 \end{aligned}$$

To bound ν from below we use that $\mathbf{E}[Z]^\dagger$ is self adjoint together with that the map $X \mapsto \langle X \mathbf{E}[Z]^\dagger X x, x \rangle$ is convex over the space of self-adjoint operators $X \in L(\mathcal{X})$ and for a fixed $x \in \mathcal{X}$. Consequently by Jensen's inequality

$$\mathbf{E} \left[\langle Z \mathbf{E}[Z]^\dagger Z x, x \rangle \right] \geq \langle \mathbf{E}[Z] \mathbf{E}[Z]^\dagger \mathbf{E}[Z] x, x \rangle \stackrel{\text{Lemma 22 item i}}{=} \langle \mathbf{E}[Z] x, x \rangle. \tag{30}$$

Finally

$$\nu \stackrel{(30)}{\geq} \sup_{x \in \text{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z] x, x \rangle}{\langle \mathbf{E}[Z] x, x \rangle} = 1.$$

Lastly, to show (14) we have

$$\begin{aligned}
\mathbf{Rank}(\mathcal{A}^*) &\stackrel{(10)}{=} \mathbf{Rank}(\mathbf{E}[Z]) \\
&\stackrel{\text{Lemma 17} + \text{Lemma 22 (v)}}{=} \mathbf{Tr}(\mathbf{E}[Z] \mathbf{E}[Z]^\dagger) = \mathbf{E}[\mathbf{Tr}(Z \mathbf{E}[Z]^\dagger)] \\
&= \mathbf{E}[\mathbf{Tr}(Z \mathbf{E}[Z]^\dagger Z)] \\
&\leq \nu \mathbf{E}[\mathbf{Tr}(Z)] \stackrel{\text{Lemma 17}}{=} \nu \mathbf{E}[\mathbf{Rank}(Z)],
\end{aligned}$$

where we used that $\langle \mathbf{E}[Z \mathbf{E}[Z]^\dagger Z] u, u \rangle \leq \nu \langle \mathbf{E}[Z] u, u \rangle$ for every $u \in \mathbf{Range}(\mathbf{E}[Z]) = \mathbf{Range}(\mathcal{A}^*) = \mathcal{X}$. \square

Proof that $X \mapsto \langle X \mathbf{E}[Z]^\dagger X x, x \rangle = \|X x\|_{\mathbf{E}[Z]^\dagger}^2$ is convex: Let $G = \mathbf{E}[Z]^\dagger$ then

$$\begin{aligned}
\|(\lambda X + (1 - \lambda)Y)x\|_G^2 &= \lambda^2 \|X x\|_G^2 + (1 - \lambda)^2 \|Y x\|_G^2 + 2\lambda(1 - \lambda) \langle x X G Y, x \rangle \\
&= -\lambda(1 - \lambda) \|(X - Y)x\|_G^2 \\
&\quad + \lambda \|X x\|_G^2 + (1 - \lambda) \|Y x\|_G^2 \\
&\leq \lambda \|X x\|_G^2 + (1 - \lambda) \|Y x\|_G^2 \quad \square.
\end{aligned}$$

A.2 Technical lemmas to prove Theorem 3

Lemma 11. For all $k \geq 0$, the vectors $y_k - x_*$, $x_k - x_*$ and $v_k - x_*$ belong to $\mathbf{Range}(\mathcal{A}^*)$.

Proof. Note that $x_0 = y_0 = x_0$ and in view of (6) we have $x_* \in x_0 + \mathbf{Range}(\mathcal{A}^*)$. So $y_0 - x_* \in \mathbf{Range}(\mathcal{A}^*)$, $v_0 - x_* \in \mathbf{Range}(\mathcal{A}^*)$ and $x_0 - x_* \in \mathbf{Range}(\mathcal{A}^*)$. Assume by induction that $y_k - x_* \in \mathbf{Range}(\mathcal{A}^*)$, $v_k - x_* \in \mathbf{Range}(\mathcal{A}^*)$ and $x_k - x_* \in \mathbf{Range}(\mathcal{A}^*)$. Since $g_k \in \mathbf{Range}(\mathcal{A}^*)$ and $x_{k+1} = y_k - g_k$ we have

$$x_{k+1} - x_* = (y_k - x_*) - g_k \in \mathbf{Range}(\mathcal{A}^*).$$

Moreover,

$$v_{k+1} - x_* = \beta(v_k - x_*) + (1 - \beta)(y_k - x_*) - \gamma g_k \in \mathbf{Range}(\mathcal{A}^*).$$

Finally

$$y_{k+1} - x_* = \alpha v_{k+1} + (1 - \alpha)x_{k+1} - x_* = \alpha(v_{k+1} - x_*) + (1 - \alpha)(x_{k+1} - x_*) \in \mathbf{Range}(\mathcal{A}^*).$$

\square

Lemma 12.

$$\mathbf{E}[\|Z_k(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2 | y_k] \leq \nu \|y_k - x_*\|_{\mathbf{E}[Z]}^2 \quad (31)$$

Proof. Since $y_k - x_* \in \mathbf{Range}(\mathcal{A}^*)$ we have that

$$\begin{aligned}
\mathbf{E}[\|Z_k(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2 | y_k] &= \langle \mathbf{E}[Z_k \mathbf{E}[Z]^\dagger Z_k] (y_k - x_*), (y_k - x_*) \rangle \\
&\stackrel{(12)}{\leq} \nu \langle \mathbf{E}[Z] (y_k - x_*), (y_k - x_*) \rangle \\
&= \nu \|y_k - x_*\|_{\mathbf{E}[Z]}^2.
\end{aligned}$$

\square

Lemma 13.

$$\|y_k - x_*\|_{\mathbf{E}[Z]}^2 = \|y_k - x_*\|^2 - \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] \quad (32)$$

Proof.

$$\begin{aligned} \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] &= \mathbf{E} [\|(I - Z_k)(y_k - x_*)\|^2 | y_k] \\ &= \langle (I - \mathbf{E}[Z])(y_k - x_*), y_k - x_* \rangle \\ &= \|y_k - x_*\|^2 - \|y_k - x_*\|_{\mathbf{E}[Z]}^2. \end{aligned}$$

□

A.3 Proof of Theorem 3

Let $r_k \stackrel{\text{def}}{=} \|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2$. It follows that

$$\begin{aligned} r_{k+1}^2 &= \|v_{k+1} - x_*\|_{\mathbf{E}[Z]^\dagger}^2 \\ &= \|\beta v_k + (1 - \beta)y_k - x_* - \gamma Z_k(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2 \\ &= \underbrace{\|\beta v_k + (1 - \beta)y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2}_I + \underbrace{\gamma^2 \|Z_k(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2}_{II} \\ &\quad - 2\gamma \underbrace{\langle \beta(v_k - x_*) + (1 - \beta)(y_k - x_*), \mathbf{E}[Z]^\dagger Z_k(y_k - x_*) \rangle}_{III} \\ &= I + \gamma^2 II - 2\gamma III. \end{aligned} \quad (33)$$

The first term can be upper bounded as follows

$$\begin{aligned} I &= \|\beta(v_k - x_*) + (1 - \beta)(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2 \\ &= \beta^2 \|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + (1 - \beta)^2 \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + 2\beta(1 - \beta) \langle v_k - x_*, y_k - x_* \rangle_{\mathbf{E}[Z]^\dagger} \\ &\stackrel{(35)}{=} \beta \|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + (1 - \beta) \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 - \beta(1 - \beta) \|v_k - y_k\|_{\mathbf{E}[Z]^\dagger}^2 \\ &\leq \beta r_k^2 + (1 - \beta) \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2, \end{aligned} \quad (34)$$

where in the third equality we used a form of the parallelogram identity

$$2\langle u, v \rangle = \|u\|^2 + \|v\|^2 - \|u - v\|^2, \quad (35)$$

with $u = v_k - x_*$ and $v = y_k - x_*$.

Taking expectation with to \mathcal{S}_k in the third term in (33) gives

$$\begin{aligned} \mathbf{E}[III | y_k, v_k, x_k] &= \langle \beta v_k + (1 - \beta)y_k - x_*, \mathbf{E}[Z]^\dagger \mathbf{E}[Z] (y_k - x_*) \rangle \\ &= \langle \beta v_k + (1 - \beta)y_k - x_*, y_k - x_* \rangle \\ &= \langle \beta \left[\frac{1}{\alpha} y_k - \frac{1 - \alpha}{\alpha} x_k \right] + (1 - \beta)y_k - x_*, y_k - x_* \rangle \\ &= \langle y_k - x_* + \beta \frac{1 - \alpha}{\alpha} (y_k - x_k), y_k - x_* \rangle \\ &= \|y_k - x_*\|^2 + \beta \frac{1 - \alpha}{\alpha} \langle y_k - x_k, y_k - x_* \rangle \\ &= \|y_k - x_*\|^2 - \beta \frac{1 - \alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_k\|^2 - \|y_k - x_*\|^2) \end{aligned} \quad (37)$$

where in the second equality (36) we used that $y_k - x_* \in \mathbf{Range}(\mathcal{A}^*) \stackrel{(10)}{=} \mathbf{Range}(\mathbf{E}[Z])$ together with a defining property of pseudoinverse operators $\mathbf{E}[Z]^\dagger \mathbf{E}[Z] w = w$ for all $w \in \mathbf{Range}(\mathbf{E}[Z])$. In the last equality (37) we used yet again the identity (35) with $u = y_k - x_k$ and $v = y_k - x_*$.

Plugging (34) and (37) into (33) and taking conditional expectation gives

$$\begin{aligned}
\mathbf{E}[r_{k+1}^2 | y_k, v_k, x_k] &= I + \gamma^2 \mathbf{E}[II | y_k] - 2\gamma \mathbf{E}[III | y_k, v_k, x_k] \\
&\stackrel{(34)+(37)+(31)}{=} \beta r_k^2 + (1 - \beta) \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \gamma^2 \nu \|y_k - x_*\|_{\mathbf{E}[Z]}^2 \\
&\quad + 2\gamma \left(-\|y_k - x_*\|^2 + \beta \frac{1 - \alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_k\|^2 - \|y_k - x_*\|^2) \right) \\
&\stackrel{(32)+(13)}{\leq} \beta r_k^2 + \frac{1 - \beta}{\mu} \|y_k - x_*\|^2 + \gamma^2 \nu (\|y_k - x_*\|^2 - \mathbf{E}[\|x_{k+1} - x_*\|^2 | y_k]) \\
&\quad + 2\gamma \left(-\|y_k - x_*\|^2 + \beta \frac{1 - \alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_*\|^2) \right). \tag{38}
\end{aligned}$$

Therefore we have that

$$\begin{aligned}
\mathbf{E}[r_{k+1}^2 + \gamma^2 \nu \|x_{k+1} - x_*\|^2 | y_k, v_k, x_k] &\leq \beta \left(r_k^2 + \underbrace{\gamma \frac{1 - \alpha}{\alpha}}_{P_1} \|x_k - x_*\|^2 \right) \\
&\quad + \underbrace{\left(\frac{1 - \beta}{\mu} - 2\gamma + \gamma^2 \nu - \beta \gamma \frac{1 - \alpha}{\alpha} \right)}_{P_2} \|y_k - x_*\|^2.
\end{aligned}$$

To establish a recurrence, we need to choose the free parameters γ, α and β so that $P_1 = \gamma^2 \nu$ and $P_2 = 0$. Furthermore we should try to set β as small as possible so as to have a fast rate of convergence. Choosing $\beta = 1 - \sqrt{\frac{\mu}{\nu}}$, $\gamma = \sqrt{\frac{1}{\mu\nu}}$, $\alpha = \frac{1}{1 + \gamma\nu}$ gives $P_2 = 0$, $\gamma^2 \nu = 1/\mu$ and

$$\mathbf{E}\left[r_{k+1}^2 + \frac{1}{\mu} \|x_{k+1} - x_*\|^2 | y_k, v_k, x_k\right] \leq \left(1 - \sqrt{\frac{\mu}{\nu}}\right) \left(r_k^2 + \frac{1}{\mu} \|x_k - x_*\|^2\right). \tag{39}$$

Taking expectation and using the tower rules gives the result. \square

A.4 Changing norm

Given an invertible positive self-adjoint $B \in L(\mathcal{X})$, suppose we want to find the least norm solution of (5) under the norm defined by $\|x\|_B \stackrel{\text{def}}{=} \sqrt{\langle Bx, x \rangle}$ as the metric in \mathcal{X} . That is, we want to solve

$$x^* \stackrel{\text{def}}{=} \arg \min_{x \in \mathcal{X}} \frac{1}{2} \|x - x_0\|_B^2, \quad \text{subject to } \mathcal{A}x = b. \tag{40}$$

By changing variables $x = B^{-1/2}z$ we have that the above is equivalent to solving

$$z^* \stackrel{\text{def}}{=} \arg \min_{z \in \mathcal{X}} \frac{1}{2} \|z - z_0\|^2, \quad \text{subject to } \mathcal{A}B^{-1/2}z = b, \tag{41}$$

with $x^* = B^{-1/2}z^*$, and $B^{1/2}$ is the unique symmetric square root of B (see Lemma 18). We can now apply Algorithm 1 to solve (41) where $\mathcal{A}B^{-1/2}$ is the system matrix. Let x_k and v_k be the resulting iterates of applying Algorithm 1. To make explicit this change in the system matrix we define the matrix

$$Z_B \stackrel{\text{def}}{=} B^{-1/2} \mathcal{A}^* \mathcal{S}_k^* (\mathcal{S}_k \mathcal{A} B^{-1} \mathcal{A}^* \mathcal{S}_k^*)^\dagger \mathcal{S}_k \mathcal{A} B^{-1/2},$$

and the constants

$$\mu_B \stackrel{\text{def}}{=} \inf_{x \in \text{Range}(B^{-1/2} \mathcal{A}^*)} \frac{\langle \mathbf{E}[Z_B] x, x \rangle}{\langle x, x \rangle} \quad (42)$$

and

$$\nu_B \stackrel{\text{def}}{=} \sup_{x \in \text{Range}(B^{-1/2} \mathcal{A}^*)} \frac{\langle \mathbf{E}[Z_B \mathbf{E}[Z_B]^\dagger Z_B] x, x \rangle}{\langle \mathbf{E}[Z_B] x, x \rangle}. \quad (43)$$

Theorem 3 then guarantees that

$$\mathbf{E} \left[\|v_{k+1} - z_*\|_{\mathbf{E}[Z_B]^\dagger}^2 + \frac{1}{\mu_B} \|x_{k+1} - z_*\|^2 \right] \leq \left(1 - \sqrt{\frac{\mu_B}{\nu_B}} \right) \mathbf{E} \left[\|v_k - z_*\|_{\mathbf{E}[Z_B]^\dagger}^2 + \frac{1}{\mu_B} \|x_k - z_*\|^2 \right].$$

Reversing our change of variables $\bar{x}_k = B^{-1/2}x_k$ and $\bar{v}_k = B^{-1/2}v_k$ in the above displayed equation gives

$$\begin{aligned} & \mathbf{E} \left[\|\bar{v}_{k+1} - x_*\|_{B^{1/2} \mathbf{E}[Z_B]^\dagger B^{1/2}}^2 + \frac{1}{\mu_B} \|\bar{x}_{k+1} - x_*\|_B^2 \right] \\ & \leq \left(1 - \sqrt{\frac{\mu_B}{\nu_B}} \right) \mathbf{E} \left[\|\bar{v}_k - x_*\|_{B^{1/2} \mathbf{E}[Z_B]^\dagger B^{1/2}}^2 + \frac{1}{\mu_B} \|\bar{x}_k - x_*\|_B^2 \right]. \end{aligned} \quad (44)$$

Thus we recover the same exact from the main theorem in [23], but in a much more general setting.

B Proof of Theorem 5

The proof follows by slight modifications of the proof of Theorem 3.

First we adapt Lemma 13. As we have $x_{k+1} - x_* = (1 - \omega Z_k)(y_k - x_*)$ the following statement follows by the same arguments as in the proof of Lemma 13.

Lemma 14 (Lemma 13').

$$\eta \|y_k - x_*\|_{\mathbf{E}[Z]}^2 = \|y_k - x_*\|^2 - \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] \quad (45)$$

Proof.

$$\begin{aligned} \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] &= \mathbf{E} [\|(I - Z_k)(y_k - x_*)\|^2 | y_k] \\ &= \mathbf{E} [\langle (I - \omega Z_k)(y_k - x_*), (I - \omega Z_k)y_k - x_* \rangle] \\ &= \|y_k - x_*\|^2 - \eta \|y_k - x_*\|_{\mathbf{E}[Z]}^2. \end{aligned}$$

□

We now follow the same steps as in proof of Theorem 3 in Section A.3. We observe, that the first time Lemma 13 is applied is in equation (38). Using Lemma 14 instead, gives

$$\begin{aligned} \mathbf{E} [r_{k+1}^2 | y_k, v_k, x_k] &\leq \beta r_k^2 + \frac{1-\beta}{\mu} \|y_k - x_*\|^2 + \frac{\gamma^2 \nu}{\eta} (\|y_k - x_*\|^2 - \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k]) \\ &\quad + 2\gamma \left(-\|y_k - x_*\|^2 + \beta \frac{1-\alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_*\|^2) \right). \end{aligned} \quad (46)$$

Therefore we have that

$$\begin{aligned} \mathbf{E} [r_{k+1}^2 + \gamma^2 \nu \|x_{k+1} - x_*\|^2 | y_k, v_k, x_k] &\leq \beta \left(r_k^2 + \underbrace{\gamma \frac{1-\alpha}{\alpha} \|x_k - x_*\|^2}_{P'_1} \right) \\ &\quad + \underbrace{\left(\frac{1-\beta}{\mu} - 2\gamma + \frac{\gamma^2 \nu}{\eta} - \beta \gamma \frac{1-\alpha}{\alpha} \right)}_{P'_2} \|y_k - x_*\|^2. \end{aligned}$$

Noting that $\frac{1-\alpha}{\alpha} = \gamma \nu$ and $\frac{\gamma^2 \nu}{\eta} = \frac{\gamma(1-\alpha)}{\eta \alpha} = \frac{1}{\mu}$, we observe $P'_2 = 0$ and deduce the statement of Theorem 5.

C Proof of Theorem 6

It suffices to study equation (38). We observe that for convergence the big bracket, P_2 , should be negative,

$$(1-\beta) \frac{1}{\mu} + \gamma^2 \nu - 2\gamma - \gamma \beta \frac{1-\alpha}{\alpha} \leq 0 \quad (47)$$

The convergence rate is then

$$\rho \stackrel{\text{def}}{=} \max \left\{ \beta, \frac{(1-\alpha)\beta}{\alpha \gamma \nu} \right\}. \quad (48)$$

or in the notation of Theorem 6, $\rho = \max\{\beta, s\beta\}$.

This means, that in order to obtain the best convergence rate, we should therefore choose parameters β and γ such that β is as small as possible. This observation is true regardless of the value of s (which itself depends on γ).

With the notation $\tau = s\gamma\beta$, we reformulate (47) to obtain

$$\frac{1}{\mu} + \gamma^2 \nu - 2\gamma \leq \beta \left(\frac{1}{\mu} + s\gamma^2 \nu \right) \quad (49)$$

Thus we see, that β cannot be chosen smaller than

$$\beta^*(s, \gamma) = \frac{1 + \mu\gamma^2 \nu - 2\mu\gamma}{1 + s\mu\gamma^2 \nu} \quad (50)$$

Minimizing this expression in γ gives

$$\beta^*(s) = \frac{1 + s - s\sqrt{\frac{\nu + 4\mu s - 2\nu s + \nu s^2}{\nu s^2}}}{2s} \quad (51)$$

with $\gamma^*(s) = \frac{1}{(1-s\beta^*(s))\nu}$.

We further observe that this parameter setting indeed guarantees convergence, i.e. $\rho \leq 1$. From (51) we observe ($\nu > 0$, $s \geq 0$, $\mu \geq 0$):

$$\beta^*(s) \leq \frac{1 + s - \sqrt{\frac{\nu - 2\nu s + \nu s^2}{\nu}}}{2s} = \frac{1 + s - (s - 1)}{2s} = \frac{1}{s} \quad (52)$$

Hence $s\beta^*(s) \leq 1$. On the other hand, $(1-s) \leq \sqrt{(1-s)^2 + \frac{4\mu s}{\nu}}$ and hence $(1+s) - \sqrt{(1-s)^2 + \frac{4\mu s}{\nu}} \leq 2s$, which shows $\beta^*(s) \leq 1$.

D Proofs and Further Comments on Section 3

D.1 Proof of Theorem 10

We perform a change of coordinates since it is easier to work with the standard Frobenius norm as opposed to the weighted Frobenius norm. Let $\hat{X} = A^{1/2}XA^{1/2}$ so that (22) and (24) become

$$\hat{X}_* \stackrel{\text{def}}{=} I = \arg \min \|\hat{X}\|_F^2 \quad \text{subject to} \quad \hat{X} = I, \quad \hat{X} = \hat{X}^\top, \quad (53)$$

and

$$\hat{X}_{k+1} = P + (I - P)\hat{X}_k(I - P), \quad (54)$$

respectively, where $P = A^{1/2}S(S^\top AS)^{-1}S^\top A^{1/2}$. The linear operator that encodes the constraint in (3.2) is given by $\hat{\mathcal{A}}(X) = (X, X - X^\top)$ the adjoint of which is given by $\hat{\mathcal{A}}^*(Y_1, Y_2) = Y_1 + Y_2 - Y_2^\top$. Since $\hat{\mathcal{A}}^*$ is clearly surjective, it follows that $\text{Range}(\hat{\mathcal{A}}^*) = \mathbb{R}^{n \times n}$.

Subtracting the identity matrix from both sides of (54) and using that P is a projection matrix, we have that

$$\hat{X}_{k+1} - I = (I - P)(\hat{X}_k - I)(I - P). \quad (55)$$

To determine the Z operator (7), from (9) and (55) we know that

$$(I - P)(\hat{X}_k - I)(I - P) = (I - Z)(\hat{X}_k - I).$$

Thus for every matrix $X \in \mathbb{R}^{n \times n}$ we have that

$$Z(X) = X - (I - P)X(I - P) = XP + PX(I - P). \quad (56)$$

Denote column-wise vectorization of X as x : $x \stackrel{\text{def}}{=} \text{Vec}(X)$. To calculate a useful lower bound on μ , note that

$$\begin{aligned} \text{Tr}(X^\top Z(X)) &= \text{Tr}(X^\top XP) + \text{Tr}(X^\top PX(I - P)) \\ &= x^\top \text{Vec}(XP) + x^\top \text{Vec}(PX(I - P)) \\ &= x^\top (P \otimes I)x + x^\top ((I - P) \otimes P)x \\ &\stackrel{(27)}{=} x^\top \mathbf{Z}x, \end{aligned} \quad (57)$$

where we used that $\text{Tr}(A^\top B) = \text{Vec}(A)^\top \text{Vec}(B)$ and $\text{Vec}(AXB) = (B^\top \otimes A)\text{Vec}(x)$ holds for any A, B, X .

Consequently, μ is equal to

$$\mu \stackrel{(11)}{=} \inf_{X \in \mathbb{R}^{n \times n}} \frac{\langle \mathbf{E}[Z] X, X \rangle_F}{\|X\|_F^2} \stackrel{(57)}{=} \inf_{x \in \mathbb{R}^{n^2 \times n^2}} \frac{x^\top \mathbf{E}[Z] x}{x^\top x} = \lambda_{\min}(\mathbf{E}[Z]).$$

Notice that we have $2\lambda_{\min}(\mathbf{E}[P]) \geq \lambda_{\min}(\mathbf{E}[Z]) \geq \lambda_{\min}(\mathbf{E}[P])$ since $(P \otimes I) + (I \otimes P) \geq Z \geq (P \otimes I)$.

In light of Algorithm 1, the iterates of the accelerated version of (54) are given by

$$\begin{aligned} \hat{Y}_k &= \alpha \hat{V}_k + (1 - \alpha) \hat{X}_k \\ \hat{G}_k &= Z_k(\hat{Y}_k - I) \\ \hat{X}_{k+1} &= \hat{Y}_k - \hat{G}_k \\ \hat{V}_{k+1} &= \beta \hat{V}_k + (1 - \beta) \hat{Y}_k - \gamma \hat{G}_k \end{aligned} \quad (58)$$

where $\hat{Y}_k, \hat{V}_k, \hat{G} \in \mathbb{R}^{n \times n}$. From Theorem 3 we have that \hat{V}_k and \hat{X}_k converge to the identity matrix according to

$$\mathbf{E} \left[\|\hat{V}_{k+1} - I\|_{\mathbf{E}[Z]^\dagger}^2 + \frac{1}{\mu} \|\hat{X}_{k+1} - I\|_F^2 \right] \leq \left(1 - \sqrt{\frac{\mu}{\nu}} \right) \mathbf{E} \left[\|\hat{V}_k - I\|_{\mathbf{E}[Z]^\dagger}^2 + \frac{1}{\mu} \|\hat{X}_k - I\|_F^2 \right], \quad (59)$$

where $\|X\|_{\mathbf{E}[Z]^\dagger}^2 = \langle \mathbf{E}[Z]^\dagger X, X \rangle_F$. Changing coordinates back to $\hat{X}_k = A^{1/2} X_k A^{1/2}$ and defining $Y_k \stackrel{\text{def}}{=} A^{-1/2} \hat{Y}_k A^{-1/2}$, $V_k \stackrel{\text{def}}{=} A^{-1/2} \hat{V}_k A^{-1/2}$ and $G_k \stackrel{\text{def}}{=} A^{-1/2} \hat{G}_k A^{-1/2}$, we have that (59) gives (25). Furthermore, using the same coordinate change applied to the iterates (58) gives Algorithm 2.

D.2 Matrix inversion as linear system

Denote $x = \text{Vec}(X)$, i.e. x is n^2 dimensional vector such that $X_{(n(i-1)+1):ni} = X_{:,i}$. Similarly, denote $e = \text{Vec}(I)$. System (4) can be thus rewritten as

$$(I \otimes A)x = e. \quad (60)$$

Notice that all linear sketches of the original system $AX = I$ can be written as

$$S_0^\top (I \otimes A)x = S_0^\top e \quad (61)$$

for a suitable $n^2 \times n^2$ matrix S_0 , therefore the setting is fairly general.

D.2.1 Alternative proof of Theorem 10

Let us now, for a purpose of this proof, consider sketch matrix S_0 to capture only sketching the original matrix system $AX = I$ by left multiplying by S , i.e. $S_0 = (I \otimes S)$, as those are the considered sketches in the setting of Section 3.

As we have

$$\text{Tr}(BX^\top BX) = \text{Vec}(BXB)^\top x = x^\top (B \otimes B)x,$$

weighted Frobenius norm of matrices is equivalent to a special weighted euclidean norm of vectors. Define also C to be a matrix such that $Cx = 0$ if and only if $X = X^\top$. Therefore, (3.2) is equivalent to

$$x_{k+1} = \arg \min \|x - x_k\|_{A \otimes A}^2 \quad \text{subject to} \quad (I \otimes S^\top)(I \otimes A)x = (I \otimes S^\top)e, \quad Cx = 0, \quad (62)$$

which is a sketch-and-project method applied on the linear system, with update as per (24):

$$x^{k+1} = x^k - (H \otimes I)((I \otimes A)x - e) - (I \otimes H)((I \otimes A)x - e) + (HA \otimes H)((I \otimes A)x - e)$$

for $H \stackrel{\text{def}}{=} S(S^\top AS)^{-1}S^\top$. Using substitution $\hat{x} = (A^{\frac{1}{2}} \otimes A^{\frac{1}{2}})x$; $\hat{S} = A^{\frac{1}{2}}S$ and comparing to (9), we get

$$Z = I \otimes I - (I - P) \otimes (I - P)$$

for P as defined inside the statement of Theorem 10. Therefore, we have all necessary information to apply the results from [23], recovering Theorem 10.

E Linear Operators in Euclidean Spaces

Here we provide some technical lemmas and results for linear operators in Euclidean space, that we used in the main body of the paper. Most of these results can be found in standard textbooks of analysis, such as [21]. We give them here for completion.

Let $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ be Euclidean spaces, equipped with inner products. Formally, we should use a notation that distinguishes the inner product in each space. But instead we use $\langle \cdot, \cdot \rangle$ to denote the inner product on all spaces, as it will be easy to determine from which space the elements are in. That is, for $x_1, x_2 \in \mathcal{X}$, we denote by $\langle x_1, x_2 \rangle$ the inner product between x_1 and x_2 in \mathcal{X} .

Let

$$\|T\| \stackrel{\text{def}}{=} \sup_{\|x\| \leq 1} \|Tx\|,$$

denote the operator norm of T . Let $0 \in L(\mathcal{X}, \mathcal{Y})$ denote the zero operator and $I \in L(\mathcal{X}, \mathcal{Y})$ the identity map.

The adjoint. Let $T^* \in L(\mathcal{Y}, \mathcal{X})$ denote the unique operator that satisfies

$$\langle Tx, y \rangle = \langle x, T^*y \rangle,$$

for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. We say that T^* is the *adjoint* of T . We say T is *self-adjoint* if $T = T^*$. Since for all $x \in \mathcal{X}$ and $s \in \mathcal{S}$,

$$\langle x, (ST)^*s \rangle = \langle STx, s \rangle_{\mathcal{S}} = \langle Tx, S^*s \rangle_{\mathcal{Y}} = \langle x, T^*S^*s \rangle,$$

we have

$$(ST)^* = T^*S^*.$$

Lemma 15. For $T \in L(\mathcal{X}, \mathcal{Y})$ we have that $\text{Range}(T^*)^\perp = \text{Null}(T)$. Thus

$$\mathcal{X} = \text{Range}(T^*) \oplus \text{Null}(T) \quad (63)$$

$$\mathcal{Y} = \text{Range}(T) \oplus \text{Null}(T^*) \quad (64)$$

Proof. See 3.2.6 in [21]. □

E.1 Positive Operators

We say that $G \in L(\mathcal{X})$ is positive if it is self-adjoint and if $\langle x, Gx \rangle \geq 0$ for all $x \in \mathcal{X}$. Let $(e_j)_{j=1}^\infty \in \mathcal{X}$ be an orthonormal basis. The trace of G is defined as

$$\mathbf{Tr}(G) \stackrel{\text{def}}{=} \sum_{j=1}^{\infty} \langle Ge_j, e_j \rangle. \quad (65)$$

The definition of trace is independent of the choice of basis due to the following lemma.

Lemma 16. *If U is unitary and $G \geq 0$ then $\mathbf{Tr}(UGU^*) = \mathbf{Tr}(G)$.*

Proof. See 3.4.3 and 3.4.4 in [21]. □

Lemma 17. *If $P \in L(\mathcal{X})$ is a projection matrix then $\mathbf{Tr}(P) = \dim(\mathbf{Range}(P)) = \mathbf{Rank}(P)$.*

Proof. Let $d = \dim(\mathbf{Range}(P))$ which is possibly infinite. Given that P is a projection we have that $\mathbf{Range}(P)$ is a closed subspace and thus there exists orthonormal basis $(e_j)_{j=1}^d$ of $\mathbf{Range}(P)$. Consequently, $\mathbf{Tr}(P) \stackrel{(65)}{=} \sum_{j=1}^d 1 = d = \dim(\mathbf{Range}(P))$. □

A square root of an operator $G \in L(\mathcal{X})$ is an operator $R \in L(\mathcal{X})$ such that $R^2 = G$.

Lemma 18. *If $G : \mathcal{X} \rightarrow \mathcal{X}$ is positive, then there exists a unique positive square root of G which we denote by $G^{1/2}$.*

Proof. See 3.2.11 in [21]. □

Lemma 19. *For any $T \in L(\mathcal{X}, \mathcal{Y})$ and any $G \in L(\mathcal{Y}, \mathcal{Y})$ that is positive and injective,*

$$\mathbf{Null}(T) = \mathbf{Null}(T^*GT), \quad (66)$$

and

$$\overline{\mathbf{Range}(T^*)} = \overline{\mathbf{Range}(T^*GT)}. \quad (67)$$

Proof. The inclusion $\mathbf{Null}(T) \subset \mathbf{Null}(T^*GT)$ is immediate. For the opposite inclusion, let $x \in \mathbf{Null}(T^*GT)$. Since G is positive we have by Lemma 18 that there exists a square root with $G^{1/2}G^{1/2} = G$. Therefore, $\langle x, T^*GTx \rangle = \langle G^{1/2}Tx, G^{1/2}Tx \rangle = 0$, which implies that $G^{1/2}Tx = 0$. Since G is injective, it follows that $G^{1/2}$ is injective and thus $x \in \mathbf{Null}(T)$. Finally (67) follows by taking the orthogonal complements of (66) and observing Lemma 15. □

As an immediate consequence of (66) and (67) we have the following lemma.

Corollary 20. *For $G : \mathcal{X} \rightarrow \mathcal{X}$ positive we have that*

$$\mathbf{Null}(G^{1/2}) = \mathbf{Null}(G) \quad (68)$$

$$\overline{\mathbf{Range}(G^{1/2})} = \overline{\mathbf{Range}(G)} \quad (69)$$

E.2 Pseudoinverse

For a bounded linear operator T define the pseudoinverse of T as follows.

Definition 21. Let $T \in L(\mathcal{X}, \mathcal{Y})$ such that $\mathbf{Range}(T)$ is closed. $T^\dagger : \mathcal{Y} \rightarrow \mathcal{X}$ is said to be the pseudoinverse if

- i) $T^\dagger T x = x$ for all $x \in \mathbf{Range}(T)$.
- ii) $T^\dagger x = 0$ for all $x \in \mathbf{Null}(T^*)$.
- iii) If $x \in \mathbf{Null}(T)$ and $y \in \mathbf{Range}(T^*)$ then $T^\dagger(x + y) = T^\dagger x + T^\dagger y$.

It follows directly from the definition (see [5] for details) that T^\dagger is a unique bounded linear operator. The following properties of pseudoinverse will be important.

Lemma 22 (Properties of pseudoinverse). Let $T \in L(\mathcal{X}, \mathcal{Y})$ such that $\mathbf{Range}(T)$ is closed. It follows that

- i) $TT^\dagger T = T$
- ii) $\mathbf{Range}(T^\dagger) = \mathbf{Range}(T^*)$ and $\mathbf{Null}(T^\dagger) = \mathbf{Null}(T^*)$
- iii) $(T^*)^\dagger = (T^\dagger)^*$
- iv) If T is self-adjoint and positive then T^\dagger is self-adjoint and positive.
- v) $T^\dagger T T^* = T^*$, that is, $T^\dagger T$ projects orthogonally onto $\mathbf{Range}(T^*)$ and along $\mathbf{Null}(T)$.
- vi) Consider the linear system $Tx = d$ where $d \in \mathbf{Range}(T)$. It follows that

$$T^\dagger d = \arg \min_{x \in \mathcal{X}} \frac{1}{2} \|x\|^2 \quad \text{subject to} \quad Tx = d. \quad (70)$$

$$\text{vii) } T^\dagger = T^*(TT^*)^\dagger$$

Proof. The proof of items *i*, *ii*, *iii*, *iv*, *v* can be found in [5]. The proof of item *vi* is alternative characterization of the pseudoinverse and it can be established by using that $d \in \mathbf{Range}(T)$ together with item *i* thus $TT^\dagger d = d$. The proof then follows by using the orthogonal decomposition $\mathbf{Range}(T^*) \oplus \mathbf{Null}(T)$ to show that $T^\dagger d$ is indeed the minimum of (70). Finally item *vii* is a direct consequence of the previous items. \square